

Fitting Mixed-Effects Models Using the lme4 Package in R

Douglas Bates

University of Wisconsin - Madison
and R Development Core Team
<Douglas.Bates@R-project.org>

International Meeting of the Psychometric Society
June 29, 2008

Outline

Organizing and plotting data; simple, scalar random effects

Mixed-modeling challenges

Models for longitudinal data

Unbalanced, non-nested data sets

Generalized linear mixed models

Evaluating the log-likelihood

Organizing data in R

- ▶ Standard rectangular data sets (columns are variables, rows are observations) are stored in *R* as *data frames*.
- ▶ The columns can be *numeric* variables (e.g. measurements or counts) or *factor* variables (categorical data) or *ordered* factor variables. These types are called the *class* of the variable.
- ▶ The `str` function provides a concise description of the structure of a data set (or any other class of object in R). The `summary` function summarizes each variable according to its class. Both are highly recommended for routine use.
- ▶ Entering just the name of the data frame causes it to be printed. For large data frames use the `head` and `tail` functions to view the first few or last few rows.

R packages

- ▶ Packages incorporate functions, data and documentation.
- ▶ You can produce packages for private or in-house use or you can contribute your package to the Comprehensive R Archive Network (CRAN), <http://cran.us.R-project.org>
- ▶ We will be using the *lme4* package from CRAN. Install it from the *Packages* menu item or with
 - > `install.packages("lme4")`
- ▶ You only need to install a package once. If a new version becomes available you can update (see the menu item).
- ▶ To use a package in an R session you attach it using
 - > `require(lme4)`
 - or
 - > `library(lme4)`(This usage causes widespread confusion of the terms “package” and “library”.)

Accessing documentation

- ▶ To be added to CRAN, a package must pass a series of quality control checks. In particular, all functions and data sets must be documented. Examples and tests can also be included.

- ▶ The `data` function provides names and brief descriptions of the data sets in a package.

```
> data(package = "lme4")
```

```
Data sets in package 'lme4':
```

```
Dyestuff           Yield of dyestuff by batch
Dyestuff2          Yield of dyestuff by batch
Pastes             Paste strength by batch and cask
Penicillin         Variation in penicillin testing
cake              Breakage angle of chocolate cakes
cbpp               Contagious bovine pleuropneumonia
sleepstudy        Reaction times in a sleep deprivation study
```

- ▶ Use `?` followed by the name of a function or data set to view its documentation. If the documentation contains an example section, you can execute it with the `example` function.

The Dyestuff data set

- ▶ The `Dyestuff`, `Penicillin` and `Pastes` data sets all come from the classic book *Statistical Methods in Research and Production*, edited by O.L. Davies and first published in 1947.
- ▶ The `Dyestuff` data are a balanced one-way classification of the `Yield` of dyestuff from samples produced from six `Batches` of an intermediate product. See `?Dyestuff`.

```
> str(Dyestuff)
```

```
'data.frame': 30 obs. of 2 variables:
 $ Batch: Factor w/ 6 levels "A","B","C","D",...: 1 1 1 1 1 2 2 2 2 2 ..
 $ Yield: num 1545 1440 1440 1520 1580 ...
```

```
> summary(Dyestuff)
```

```
Batch      Yield
A:5   Min.    :1440
B:5   1st Qu.:1469
C:5   Median  :1530
D:5   Mean    :1528
E:5   3rd Qu.:1575
F:5   Max.    :1635
```

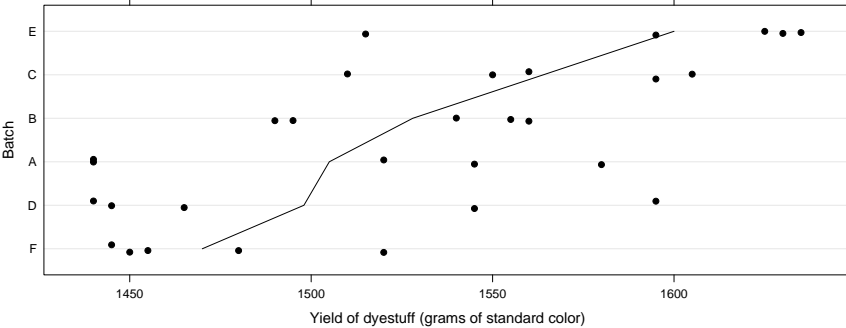
Lattice graphics

- ▶ One of the strengths of R is its graphics capabilities.
- ▶ There are several styles of graphics in R. The style in Deepayan Sarkar's *lattice* package is well-suited to the type of data we will be discussing.
- ▶ I will not show every piece of code used to produce the data graphics. The code is available in the script files for the slides (and sometimes in the example sections of the data set's documentation).
- ▶ Deepayan's book, *Lattice: Multivariate Data Visualization with R* (Springer, 2008) provides in-depth documentation and explanations of lattice graphics.
- ▶ I also recommend Phil Spector's book, *Data Manipulation with R* (Springer, 2008).

The effect of the batches

- ▶ To emphasize that `Batch` is categorical, we use letters instead of numbers to designate the levels.
- ▶ Because there is no inherent ordering of the levels of `Batch`, we will reorder the levels if, say, doing so can make a plot more informative.
- ▶ The particular batches observed are just a selection of the possible batches and are entirely used up during the course of the experiment.
- ▶ It is not particularly important to estimate and compare yields from these batches. Instead we wish to estimate the variability in yields due to batch-to-batch variability.
- ▶ The `Batch` factor will be used in *random-effects* terms in models that we fit.

Dyestuff data plot



- ▶ The line joins the mean yields of the six batches, which have been reordered by increasing mean yield.
- ▶ The vertical positions are jittered slightly to reduce overplotting. The lowest yield for batch A was observed on two distinct preparations from that batch.

Extracting information from the fitted model

- ▶ `fm1` is an object of class "mer" (mixed-effects representation).
- ▶ There are many *extractor* functions that can be applied to such objects.

```
> fixef(fm1)
(Intercept)
1527.5

> ranef(fm1, drop = TRUE)
$Batch
      A      B      C      D      E      F
-17.60597  0.39124 28.56079 -23.08338 56.73033 -44.99302

> fitted(fm1)
 [1] 1509.9 1509.9 1509.9 1509.9 1509.9 1527.9 1527.9 1527.9
 [9] 1527.9 1527.9 1556.1 1556.1 1556.1 1556.1 1556.1 1504.4
[17] 1504.4 1504.4 1504.4 1504.4 1584.2 1584.2 1584.2 1584.2
[25] 1584.2 1482.5 1482.5 1482.5 1482.5 1482.5
```

A mixed-effects model for the dyestuff yield

```
> fm1 <- lmer(Yield ~ 1 + (1 | Batch), Dyestuff)
> print(fm1)
```

```
Linear mixed model fit by REML
Formula: Yield ~ 1 + (1 | Batch)
Data: Dyestuff
   AIC   BIC logLik deviance REMLdev
325.7 329.9 -159.8   327.4   319.7

Random effects:
Groups   Name      Variance Std.Dev.
Batch   (Intercept) 1763.7   41.996
Residual                2451.3   49.511

Number of obs: 30, groups: Batch, 6

Fixed effects:
              Estimate Std. Error t value
(Intercept) 1527.50      19.38    78.81
```

- ▶ Fitted model `fm1` has one fixed-effect parameter, the mean yield, and one random-effects term, generating a simple, scalar random effect for each level of `Batch`.

Definition of linear mixed-effects models

- ▶ A mixed-effects model incorporates two vector-valued random variables: the response, \mathcal{Y} , and the random effects, \mathcal{B} . We observe the value, y , of \mathcal{Y} . We do not observe the value of \mathcal{B} .
- ▶ In a *linear mixed-effects model* the conditional distribution, $\mathcal{Y}|\mathcal{B}$, and the marginal distribution, \mathcal{B} , are independent, multivariate normal (or "Gaussian") distributions,

$$(\mathcal{Y}|\mathcal{B} = b) \sim \mathcal{N}(\mathbf{X}\beta + \mathbf{Z}b, \sigma^2\mathbf{I}), \quad \mathcal{B} \sim \mathcal{N}(\mathbf{0}, \sigma^2\Sigma), \quad (\mathcal{Y}|\mathcal{B}) \perp \mathcal{B}.$$

- ▶ The scalar σ is the *common scale parameter*; the p -dimensional β is the *fixed-effects parameter*; the $n \times p$ \mathbf{X} and the $n \times q$ \mathbf{Z} are known, fixed *model matrices*; and the $q \times q$ *relative variance-covariance matrix* $\Sigma(\theta)$ is a positive semidefinite, symmetric $q \times q$ matrix that depends on the parameter θ .

Formulation of the marginal variance matrix

- ▶ In addition to determining Z , the random effects terms determine the form and parameterization of the relative variance-covariance matrix, $\Sigma(\theta)$.
- ▶ The parameterization is based on a modified "LDL" Cholesky factorization

$$\Sigma = TSS'T'$$

where T is a $q \times q$ unit lower triangular matrix and S is a $q \times q$ diagonal Scale matrix with nonnegative diagonal elements.

- ▶ Σ , T and S are all block-diagonal, with blocks corresponding to the random-effects terms.
- ▶ The diagonal block of T for a scalar random effects term is the identity matrix, I , and the block in S is a nonnegative multiple of I .

Verbose fitting, extracting T and S

- ▶ The optional argument `verbose = TRUE` causes `lmer` to print iteration information during the optimization of the parameter estimates.
- ▶ The quantity being minimized is the *profiled deviance* of the model. The deviance is negative twice the log-likelihood. It is profiled in the sense that it is a function of θ only — β and σ are at their conditional estimates.
- ▶ If you want to see exactly how the parameters θ generate Σ , use `expand` to obtain a list with components `sigma`, `T` and `S`. The list also contains a permutation matrix `P` whose role we will discuss later.
- ▶ T , S and Σ can be very large but are always highly patterned. The `image` function can be used to examine their structure.

Obtain the verbose output for fitting fm1

```
> invisible(update(fm1, verbose = TRUE))
0:      319.76562: 0.730297
1:      319.73553: 0.962418
2:      319.65736: 0.869480
3:      319.65441: 0.844020
4:      319.65428: 0.848469
5:      319.65428: 0.848327
6:      319.65428: 0.848324
```

- ▶ The first number on each line is the iteration count — iteration 0 is at the starting value for θ .
- ▶ The second number is the profiled deviance — the criterion to be minimized at the estimates.
- ▶ The third and subsequent numbers are the parameter vector θ .

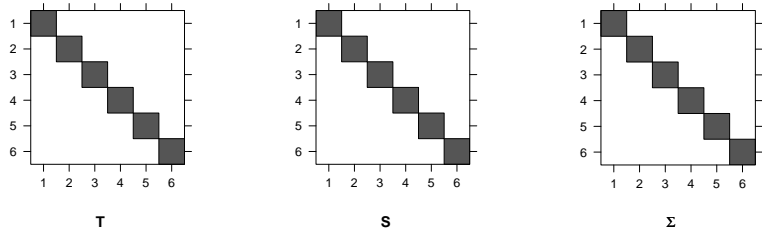
Extract T and S

```
▶ As previously indicated, T and S from fm1 are boring.
> efm1 <- expand(fm1)
> efm1$S
6 x 6 diagonal matrix of class "ddiMatrix"
      [,1] [,2] [,3] [,4] [,5] [,6]
[1,] 0.84823      .      .      .      .      .
[2,]      . 0.84823      .      .      .      .
[3,]      .      . 0.84823      .      .      .
[4,]      .      .      . 0.84823      .      .
[5,]      .      .      .      . 0.84823      .
[6,]      .      .      .      .      . 0.84823

> efm1$T
6 x 6 sparse Matrix of class "dtCMatrix"
[1,] 1 . . . . .
[2,] . 1 . . . .
[3,] . . 1 . . .
[4,] . . . 1 . .
[5,] . . . . 1 .
[6,] . . . . . 1
```

Reconstructing Σ

```
> (fm1S <- tcrossprod(efm1$T %*% efm1$S))
6 x 6 sparse Matrix of class "dsCMatrix"
[1,] 0.71949 . . . . .
[2,] . 0.71949 . . . . .
[3,] . . 0.71949 . . . . .
[4,] . . . 0.71949 . . . . .
[5,] . . . . 0.71949 . . . . .
[6,] . . . . . 0.71949 . . . . .
```



REML estimates versus ML estimates

- ▶ The default parameter estimation criterion for linear mixed models is restricted (or “residual”) maximum likelihood (REML).
- ▶ Maximum likelihood (ML) estimates (sometimes called “full maximum likelihood”) can be requested by specifying `REML = FALSE` in the call to `lmer`.
- ▶ Generally REML estimates of variance components are preferred. ML estimates are known to be biased. Although REML estimates are not guaranteed to be unbiased, they are usually less biased than ML estimates.
- ▶ Roughly the difference between REML and ML estimates of variance components is comparable to estimating σ^2 in a fixed-effects regression by $SSR/(n - p)$ versus SSR/n , where SSR is the residual sum of squares.
- ▶ For a balanced, one-way classification like the `Dyestuff` data, the REML and ML estimates of the fixed-effects are identical.

Re-fitting the model for ML estimates

```
> (fm1M <- update(fm1, REML = FALSE))
Linear mixed model fit by maximum likelihood
Formula: Yield ~ 1 + (1 | Batch)
Data: Dyestuff
   AIC   BIC logLik deviance REMLdev
333.3 337.5 -163.7   327.3   319.7
Random effects:
Groups   Name      Variance Std.Dev.
Batch    (Intercept) 1388.1   37.258
Residual                    2451.3   49.511
Number of obs: 30, groups: Batch, 6
Fixed effects:
      Estimate Std. Error t value
(Intercept) 1527.50    17.69   86.33
(The extra parentheses around the assignment cause the value to be printed. Generally the results of assignments are not printed.)
```

Recap of the Dyestuff model

- ▶ The model is fit as `lmer(formula = Yield ~ 1 + (1 | Batch), data = Dyestuff)`
- ▶ There is one random-effects term, `(1|Batch)`, in the model formula. It is a simple, scalar term for the grouping factor `Batch` with $n_1 = 6$ levels. Thus $q = 6$.
- ▶ The model matrix Z is the 30×6 matrix of indicators of the levels of `Batch`.
- ▶ The relative variance-covariance matrix, Σ , is a nonnegative multiple of the 6×6 identity matrix I_6 .
- ▶ The fixed-effects parameter vector, β , is of length $p = 1$. All the elements of the 30×1 model matrix X are unity.

The Penicillin data (also check the ?Penicillin description)

Penicillin data plot

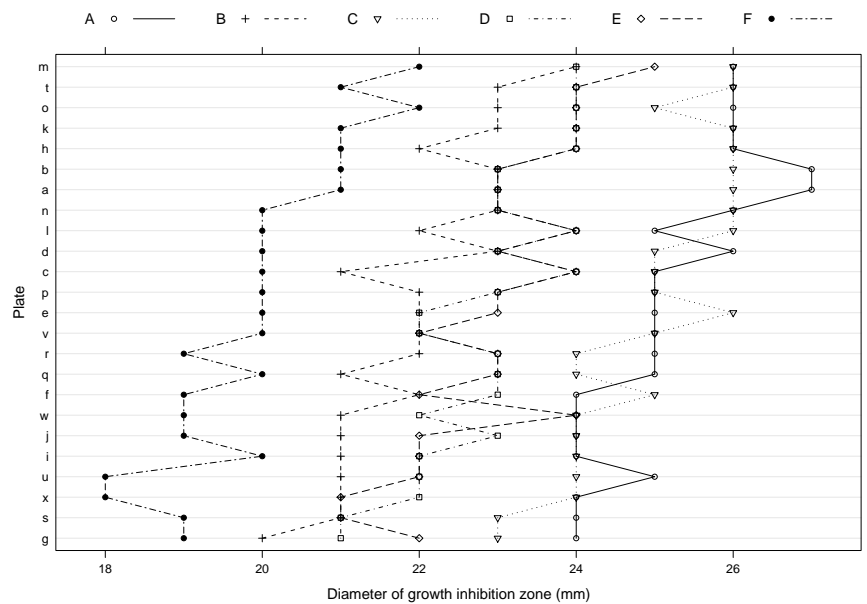
```
> str(Penicillin)

'data.frame': 144 obs. of 3 variables:
 $ diameter: num 27 23 26 23 23 21 27 23 26 23 ...
 $ plate : Factor w/ 24 levels "a","b","c","d",...: 1 1 1 1 1 1 2 2 2
 $ sample : Factor w/ 6 levels "A","B","C","D",...: 1 2 3 4 5 6 1 2 3 4

> xtabs(~sample + plate, Penicillin)
```

sample	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o	p	q	r	s	t	u	v	w	x
A	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
B	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
C	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
D	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
E	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
F	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1

- ▶ These are measurements of the potency (measured by the diameter of a clear area on a Petri dish) of penicillin samples in a balanced, unreplicated two-way crossed classification with the test medium, plate.



Model with crossed simple random effects for Penicillin

```
> (fm2 <- lmer(diameter ~ 1 + (1 | plate) + (1 | sample),
+ Penicillin))
```

```
Linear mixed model fit by REML
Formula: diameter ~ 1 + (1 | plate) + (1 | sample)
Data: Penicillin
AIC BIC logLik deviance REMLdev
338.9 350.7 -165.4 332.3 330.9

Random effects:
Groups Name Variance Std.Dev.
plate (Intercept) 0.71691 0.84671
sample (Intercept) 3.73030 1.93140
Residual 0.30242 0.54992

Number of obs: 144, groups: plate, 24; sample, 6
Fixed effects:
Estimate Std. Error t value
(Intercept) 22.9722 0.8085 28.41
```

Fixed and random effects for fm2

- ▶ The model for the $n = 144$ observations has $p = 1$ fixed-effects parameter and $q = 30$ random effects from $k = 2$ random effects terms in the formula.

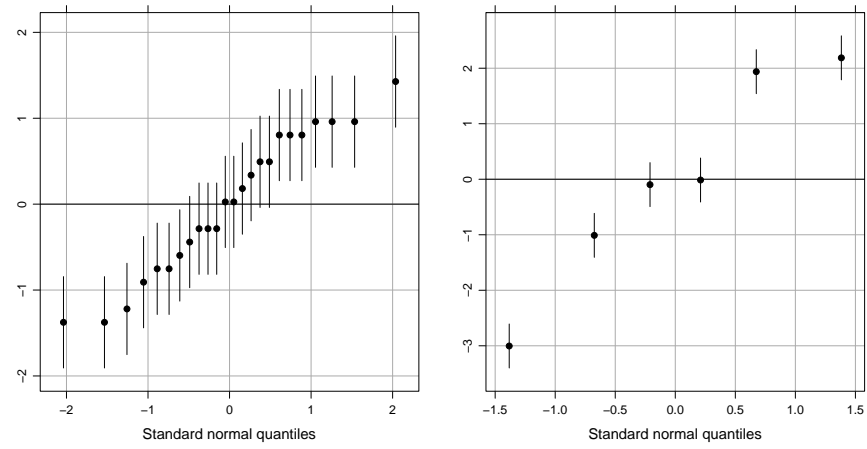
```
> fixef(fm2)
(Intercept)
22.972

> ranef(fm2, drop = TRUE)

$plate
a b c d e f
0.804547 0.804547 0.181672 0.337391 0.025953 -0.441203
g h i j k l
-1.375516 0.804547 -0.752641 -0.752641 0.960266 0.493109
m n o p q r
1.427422 0.493109 0.960266 0.025953 -0.285484 -0.285484
s t u v w x
-1.375516 0.960266 -0.908360 -0.285484 -0.596922 -1.219797

$sample
A B C D E F
2.187057 -1.010476 1.937898 -0.096895 -0.013842 -3.003742
```

Prediction intervals for random effects



Models with crossed random effects

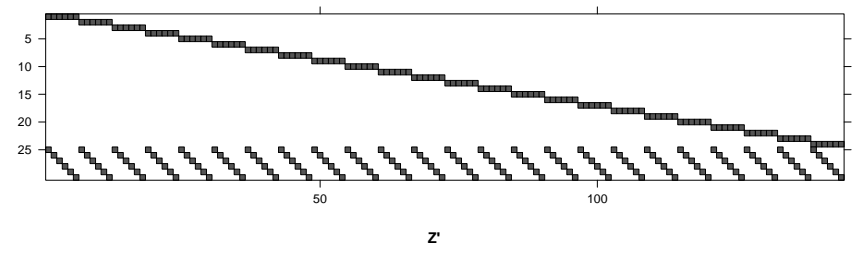
- ▶ Many people believe that mixed-effects models are equivalent to hierarchical linear models (HLMs) or “multilevel models”. This is not true. The `plate` and `sample` factors in `fm2` are crossed. They do not represent levels in a hierarchy.
- ▶ There is no difficulty in defining and fitting models with crossed random effects (meaning random-effects terms whose grouping factors are crossed).
- ▶ Crossing of random effects can affect the speed with which a model can be fit.
- ▶ The crucial calculation in each `lmer` iteration is evaluation of the sparse, lower triangular, Cholesky factor, $L(\theta)$, that satisfies

$$L(\theta)L(\theta)' = P(A(\theta)A(\theta)' + I_q)P'$$

from $A(\theta)' = ZT(\theta)S(\theta)$. Crossing of grouping factors increases the number of nonzeros in AA' and also causes some “fill-in” when creating L from A .

Model matrix Z for fm2

- ▶ Because the model matrix Z is generated from $k = 2$ simple, scalar random effects terms, it consists of two sets of indicator columns.
- ▶ The structure of Z' is shown below. (Generally we will show the transpose of these model matrices - they fit better on slides.)

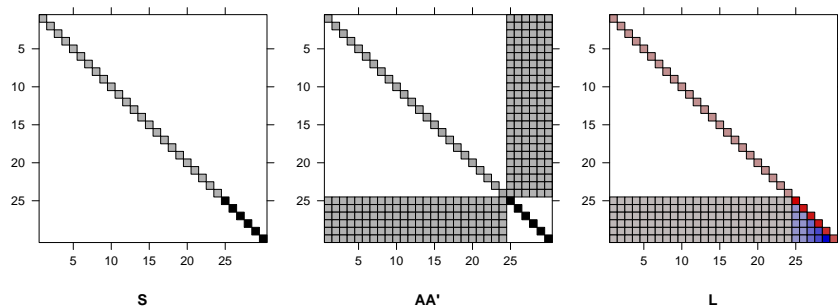


All HLMs are mixed models but not vice-versa

- ▶ Even though Raudenbush and Bryk do discuss models for crossed factors in their HLM book, such models are not hierarchical.
- ▶ Experimental situations with crossed random factors, such as “subject” and “stimulus”, are common. We can and should model such data according to its structure.
- ▶ In longitudinal studies of subjects in social contexts (e.g. students in classrooms or in schools) we almost always have partial crossing of the subject and the context factors, meaning that, over the course of the study, a particular student may be observed in more than one class (partial crossing) but not all students are observed in all classes. The student and class factors are neither fully crossed nor strictly nested.
- ▶ For longitudinal data, “nested” is only important if it means “nested across time”. “Nested at a particular time” doesn’t count.
- ▶ The `lme4` package in `R` is different from most other software

Images of some of the $q \times q$ matrices for fm2

- ▶ Because both random-effects terms are scalar terms, \mathbf{T} is a block-diagonal matrix of two blocks, both of which are identity matrices. Hence $\mathbf{T} = \mathbf{I}_q$.
- ▶ For this model it is also the case that $\mathbf{P} = \mathbf{I}_q$.
- ▶ \mathbf{S} consists of two diagonal blocks, both of which are multiples of an identity matrix. The multiples are different.



Recap of the Penicillin model

- ▶ The model formula is $\text{diameter} \sim 1 + (1 | \text{plate}) + (1 | \text{sample})$
- ▶ There are two random-effects terms, $(1 | \text{plate})$ and $(1 | \text{sample})$. Both are simple, scalar ($q_1 = q_2 = 1$) random effects terms, with $n_1 = 24$ and $n_2 = 6$ levels, respectively. Thus $q = q_1 n_1 + q_2 n_2 = 30$.
- ▶ The model matrix \mathbf{Z} is the 144×30 matrix created from two sets of indicator columns.
- ▶ The relative variance-covariance matrix, Σ , is block diagonal in two blocks that are nonnegative multiples of identity matrices. The matrices \mathbf{AA}' and \mathbf{L} show the crossing of the factors. \mathbf{L} has some fill-in relative to \mathbf{AA}' .
- ▶ The fixed-effects parameter vector, β , is of length $p = 1$. All the elements of the 144×1 model matrix \mathbf{X} are unity.

The Pastes data (also check the ?Pastes description)

```
> str(Pastes)
```

```
'data.frame': 60 obs. of 4 variables:
 $ strength: num  62.8 62.6 60.1 62.3 62.7 63.1 60 61.4 57.5 56.9 ...
 $ batch   : Factor w/ 10 levels "A","B","C","D",...: 1 1 1 1 1 1 2 2 2
 $ cask    : Factor w/ 3 levels "a","b","c": 1 1 2 2 3 3 1 1 2 2 ...
 $ sample  : Factor w/ 30 levels "A:a","A:b","A:c",...: 11 2 2 3 3 4 4
```

```
> xtabs(~batch + sample, Pastes, sparse = TRUE)
```

```
10 x 30 sparse Matrix of class "dgCMatrix"
A 2 2 2 . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
B . . . 2 2 2 . . . . . . . . . . . . . . . . . . . . . . . . . . . .
C . . . . . 2 2 2 . . . . . . . . . . . . . . . . . . . . . . . . . .
D . . . . . . 2 2 2 . . . . . . . . . . . . . . . . . . . . . . . . . .
E . . . . . . . 2 2 2 . . . . . . . . . . . . . . . . . . . . . . . .
F . . . . . . . . 2 2 2 . . . . . . . . . . . . . . . . . . . . . . . .
G . . . . . . . . . 2 2 2 . . . . . . . . . . . . . . . . . . . . . . .
H . . . . . . . . . . 2 2 2 . . . . . . . . . . . . . . . . . . . . . .
I . . . . . . . . . . . 2 2 2 . . . . . . . . . . . . . . . . . . . . .
J . . . . . . . . . . . . 2 2 2 . . . . . . . . . . . . . . . . . . . .
```

Structure of the Pastes data

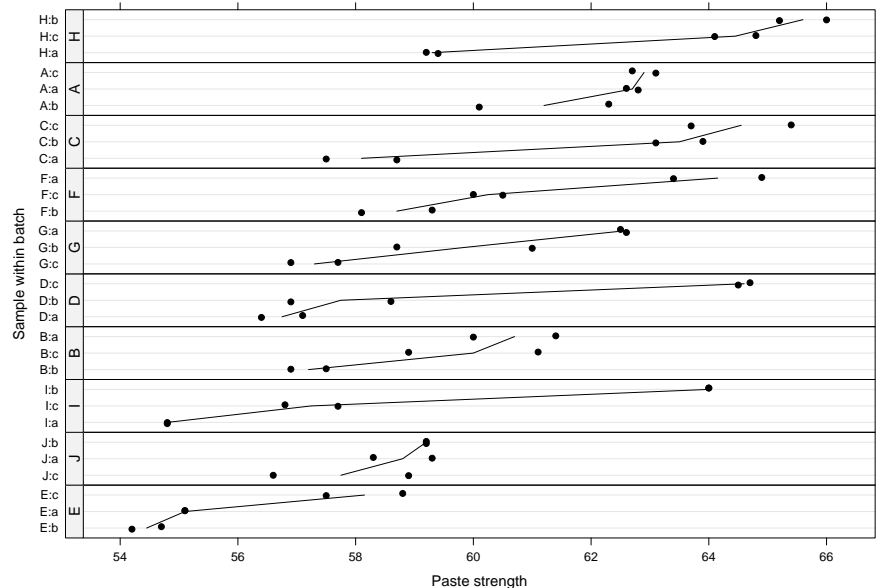
- ▶ The `sample` factor is nested within the `batch` factor. Each sample is from one of three casks selected from a particular batch.
- ▶ Note that there are 30, not 3, distinct samples.
- ▶ We can label the casks as 'a', 'b' and 'c' but then the `cask` factor by itself is meaningless (because cask 'a' in batch 'A' is unrelated to cask 'a' in batches 'B', 'C', ...). The `cask` factor is only meaningful within a `batch`.
- ▶ Only the `batch` and `cask` factors, which are apparently crossed, were present in the original data set. `cask` may be described as being nested within `batch` but that is not reflected in the data. It is *implicitly nested*, not explicitly nested.
- ▶ You can save yourself a lot of grief by immediately creating the explicitly nested factor. The recipe is

```
> Pastes <- within(Pastes, sample <- (batch:cask)[drop = TRUE])
```

Avoid implicitly nested representations

- ▶ The `lme4` package allows for very general model specifications. It does not require that factors associated with random effects be hierarchical or “multilevel” factors in the design.
- ▶ The same model specification can be used for data with nested or crossed or partially crossed factors. Nesting or crossing is determined from the structure of the factors in the data, not the model specification.
- ▶ You can avoid confusion about nested and crossed factors by following one simple rule: ensure that different levels of a factor in the experiment correspond to different labels of the factor in the data.
- ▶ Samples were drawn from 30, not 3, distinct casks in this experiment. We should specify models using the `sample` factor with 30 levels, not the `cask` factor with 3 levels.

Pastes data plot



A model with nested random effects

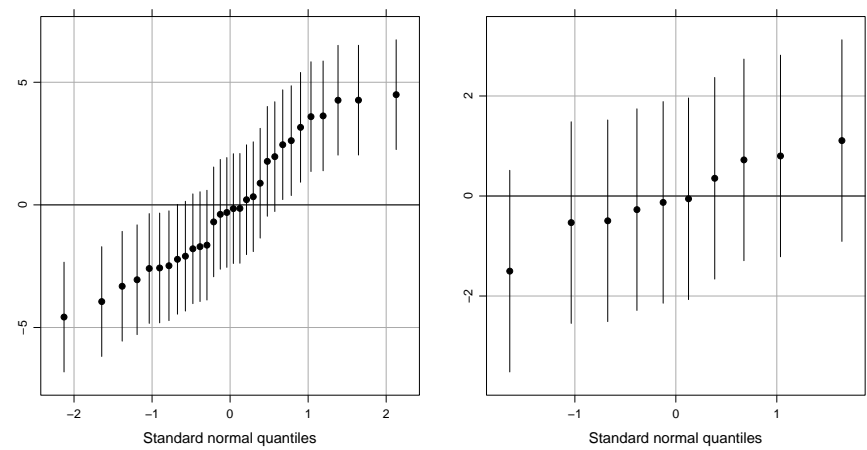
```
> (fm3 <- lmer(strength ~ 1 + (1 | batch) + (1 | sample),
+ Pastes))
```

```
Linear mixed model fit by REML
Formula: strength ~ 1 + (1 | batch) + (1 | sample)
Data: Pastes
AIC   BIC logLik deviance REMLdev
255 263.4 -123.5  248.0    247

Random effects:
Groups   Name      Variance Std.Dev.
sample  (Intercept) 8.43378  2.90410
batch   (Intercept) 1.65691  1.28721
Residual                    0.67801  0.82341

Number of obs: 60, groups: sample, 30; batch, 10
Fixed effects:
              Estimate Std. Error t value
(Intercept)  60.0533    0.6768   88.73
```

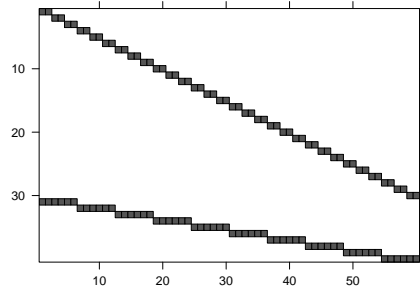
Random effects from model fm3



- ▶ This plot and the data plot both show that the sample-to-sample variability dominates the batch-to-batch variability. We will return to this observation later.

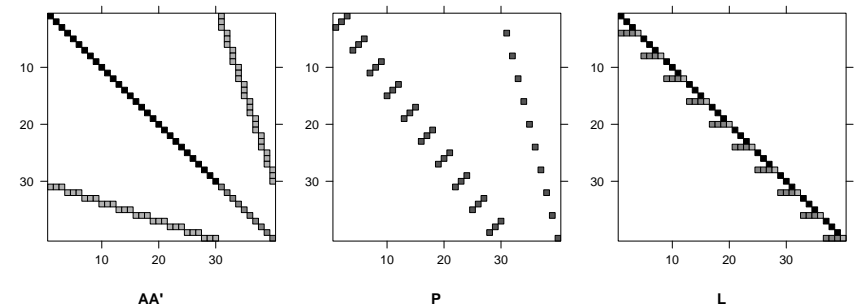
Dimensions and relationships in fm3

- ▶ There are $n = 60$ observations, $p = 1$ fixed-effects parameter, $k = 2$ simple, scalar random-effects terms ($q_1 = q_2 = 1$) with grouping factors having $n_1 = 30$ and $n_2 = 10$ levels.
- ▶ Because both random-effects terms are scalar terms, $\mathbf{T} = \mathbf{I}_{40}$ and \mathbf{S} is block-diagonal in two diagonal blocks of sizes 30 and 10, respectively. \mathbf{Z} is generated from two sets of indicators.



Images of some of the $q \times q$ matrices for fm3

- ▶ The permutation \mathbf{P} has two purposes: reduce fill-in and “post-order” the columns to keep nonzeros near the diagonal.
- ▶ In a model with strictly nested grouping factors there will be no fill-in. The permutation \mathbf{P} is chosen for post-ordering only.



Eliminate the random-effects term for batch?

- ▶ We have seen that there is little batch-to-batch variability beyond that induced by the variability of samples within batches.
- ▶ We can fit a reduced model without that term and compare it to the original model.
- ▶ Somewhat confusingly, model comparisons from likelihood ratio tests are obtained by calling the `anova` function on the two models. (Put the simpler model first in the call to `anova`.)
- ▶ Sometimes likelihood ratio tests can be evaluated using the REML criterion and sometimes they can't. Instead of learning the rules of when you can and when you can't, it is easiest always to refit the models with `REML = FALSE` before comparing.

Comparing ML fits of the full and reduced models

```
> fm3M <- update(fm3, REML = FALSE)
> fm4M <- lmer(strength ~ 1 + (1 | sample), Pastes,
+             REML = FALSE)
> anova(fm4M, fm3M)
```

Data: Pastes

Models:

fm4M: strength ~ 1 + (1 | sample)

fm3M: strength ~ 1 + (1 | batch) + (1 | sample)

	Df	AIC	BIC	logLik	Chisq	Chi	Df	Pr(>Chisq)
fm4M	3	254.40	260.69	-124.20				
fm3M	4	255.99	264.37	-124.00	0.4072		1	0.5234

p-values of LR tests on variance components

- ▶ The likelihood ratio is a reasonable criterion for comparing these two models. However, the theory behind using a χ^2 distribution with 1 degree of freedom as a reference distribution for this test statistic does not apply in this case. The null hypothesis is on the boundary of the alternative hypothesis.
- ▶ Even at the best of times, the p-values for such tests are only approximate because they are based on the asymptotic behavior of the test statistic. To carry the argument further, all results in statistics are based on models and, as George Box famously said, “All models are wrong; some models are useful.”

LR tests on variance components (cont'd)

- ▶ In this case the problem with the boundary condition results in a p-value that is larger than it would be if, say, you compared this likelihood ratio to values obtained for data simulated from the null hypothesis model. We say these results are “conservative”.
- ▶ As a rule of thumb, the p-value for a simple, scalar term is roughly twice as large as it should be.
- ▶ In this case, dividing the p-value in half would not affect our conclusion.

Updated model, REML estimates

```
> (fm4 <- update(fm4M, REML = TRUE))
```

```
Linear mixed model fit by REML
Formula: strength ~ 1 + (1 | sample)
Data: Pastes
   AIC   BIC logLik deviance REMLdev
253.6 259.9 -123.8   248.4   247.6
Random effects:
Groups   Name             Variance Std.Dev.
sample  (Intercept)  9.97622  3.15852
Residual                    0.67803  0.82342
Number of obs: 60, groups: sample, 30
Fixed effects:
              Estimate Std. Error t value
(Intercept)  60.0533    0.5864   102.4
```

Recap of the analysis of the Pastes data

- ▶ The data consist of $n = 60$ observations on $q_1 = 30$ samples nested within $q_2 = 10$ batches.
- ▶ The data are labelled with a `batch` factor with 3 levels but that is an implicitly nested factor. Create the explicit factor `sample` and ignore `batch` from then on.
- ▶ Specification of a model for nested factors is exactly the same as specification of a model with crossed or partially crossed factors — provided that you avoid using implicitly nested factors.
- ▶ In this case the `batch` factor was inert — it did not “explain” substantial variability in addition to that attributed to the `sample` factor. We therefore prefer the simpler model.
- ▶ At the risk of “beating a dead horse”, notice that, if we had used the `batch` factor in some way, we would still need to create a factor like `sample` to be able to reduce the model. The `batch` factor is only meaningful within `batch`.

Recap of simple, scalar random-effects terms

- ▶ For the `lmer` function (and also for `glmer` and `nlmer`) a simple, scalar random effects term is of the form $(1|F)$.
- ▶ The number of random effects generated by the i th such term is the number of levels, n_i , of F (after dropping “unused” levels — those that do not occur in the data. The idea of having such levels is not as peculiar as it may seem if, say, you fitting a model to a subset of the original data.)
- ▶ Such a term contributes n_i columns to Z . These columns are the indicator columns of the grouping factor.
- ▶ Such a term contributes a diagonal block I_{n_i} to T . If all random effects terms are scalar terms then $T = I$.
- ▶ Such a term contributes a diagonal block $c_i I_{n_i}$ to S . The multipliers c_i can be different for different terms. The term contributes exactly one element (which is c_i) to θ .

A large observational data set

- ▶ A major university (not mine) provided data on the grade point score (`gr.pt`) by student (`id`), instructor (`instr`) and department (`dept`) from a 10 year period. I regret that I cannot make these data available to others.
- ▶ These factors are unbalanced and partially crossed.

```
> str(anon.grades.df)
```

```
'data.frame': 1721024 obs. of 9 variables:
 $ instr  : Factor w/ 7964 levels "10000","10001",...: 1 1 1 1 1 1 1 1 1
 $ dept   : Factor w/ 106 levels "AERO","AFAM",...: 43 43 43 43 43 43 4 4
 $ id     : Factor w/ 54711 levels "900000001","900000002",...: 12152 1
 $ nclass : num  40 29 33 13 47 49 37 14 21 20 ...
 $ vgpa   : num  NA NA NA NA NA NA NA NA NA NA ...
 $ rawai  : num  2.88 -1.15 -0.08 -1.94 3.00 ...
 $ gr.pt  : num  4 1.7 2 0 3.7 1.7 2 4 2 2.7 ...
 $ section: Factor w/ 70366 levels "19959 AERO011A001",...: 18417 18417
 $ semester: num  19989 19989 19989 19989 19972 ...
```

This is all very nice, but ...

- ▶ These methods are interesting but the results are not really new. Similar results are quoted in *Statistical Methods in Research and Production*, which is a very old book.
- ▶ The approach described in that book is actually quite sophisticated, especially when you consider that the methods described there, based on observed and expected mean squares, are for hand calculation (in pre-calculator days)!
- ▶ Why go to all the trouble of working with sparse matrices and all that if you could get the same results with paper and pencil? The one-word answer is *balance*.
- ▶ Those methods depend on the data being balanced. The design must be completely balanced and the resulting data must also be completely balanced.
- ▶ Balance is fragile. Even if the design is balanced, a single missing or questionable observation destroys the balance. Observational studies (as opposed to, say, laboratory experiments) cannot be expected to yield balanced data sets.
- ▶ Also, the models involve only simple, scalar random effects

A preliminary model

```
Linear mixed model fit by REML
Formula: gr.pt ~ (1 | id) + (1 | instr) + (1 | dept)
Data: anon.grades.df
      AIC      BIC    logLik deviance REMLdev
3447389 3447451 -1723690  3447374 3447379
Random effects:
Groups   Name             Variance Std.Dev.
id       (Intercept)  0.3085  0.555
instr    (Intercept)  0.0795  0.282
dept     (Intercept)  0.0909  0.301
Residual                    0.4037  0.635
Number of obs: 1685394, groups: id, 54711; instr, 7915; dept, 102

Fixed effects:
              Estimate Std. Error t value
(Intercept)   3.1996      0.0314    102
```

Comments on the model fit

- ▶ $n = 1685394$, $p = 1$, $k = 3$, $n_1 = 54711$, $n_2 = 7915$, $n_3 = 102$, $q_1 = q_2 = q_3 = 1$, $q = 62728$
- ▶ This model is sometimes called the “unconditional” model in that it does not incorporate covariates beyond the grouping factors.
- ▶ It takes less than an hour to fit an “unconditional” model with random effects for student (`id`), instructor (`inst`) and department (`dept`) to these data.
- ▶ Naturally, this is just the first step. We want to look at possible time trends and the possible influences of the covariates.

Looking at the “big picture”

- ▶ Most descriptions of mixed-effects models concentrate on the details, sometimes degenerating into subscript fests. As in many complex systems, it is valuable before considering the details to step back and look at the big picture.
- ▶ We consider mixed models for an n -dimensional response \mathbf{y} modeled as a random variable \mathcal{Y} . The random effects, \mathcal{B} , are another vector-valued random variable in the model.
- ▶ The model specifies the conditional distribution, $\mathcal{Y}|\mathcal{B}$, and the marginal distribution, \mathcal{B} , as depending on parameters.
- ▶ The marginal distribution, \mathcal{B} , has the form

$$\mathcal{B} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \Sigma(\boldsymbol{\theta}))$$

where $\boldsymbol{\theta}$ is the *variance-component* parameter vector. Not all models incorporate a *common scale parameter*, σ . When it is used, it also occurs in the formulation of the conditional distribution, $\mathcal{Y}|\mathcal{B}$.

Challenges in fitting mixed models

- ▶ Like all statistical models, linear mixed models are being fit to larger and larger data sets - microarray data, massive longitudinal studies, genetic studies with pedigrees, etc.
- ▶ The structure of the models being fit is also getting more complex. Observational longitudinal studies with multiple levels of grouping (e.g. test scores by student, teacher, school, district, ...) may have non-nested groupings (a particular student is exposed to more than one teacher/school combination). Or we could be modeling data classified by fully crossed factors such as subject and item.
- ▶ We wish to allow for generalizations such as generalized linear mixed models (GLMMs) or nonlinear mixed models (NLMMs) or even generalized nonlinear mixed models (GNMMs).

The conditional distribution, $\mathcal{Y}|\mathcal{B}$

- ▶ The conditional mean,

$$\mu_{\mathcal{Y}|\mathcal{B}}(\mathbf{b}) = \mathbb{E}[\mathcal{Y}|\mathcal{B} = \mathbf{b}] = \mathbf{g}^{-1}(\boldsymbol{\eta}),$$

depends on the *unbounded predictor*, $\boldsymbol{\eta}$, through the *inverse link*, \mathbf{g}^{-1} . For LMMs and GLMMs, $\boldsymbol{\eta}$ is a linear predictor,

$$\boldsymbol{\eta} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{b}.$$

- ▶ The conditional distribution, $(\mathcal{Y}|\mathcal{B} = \mathbf{b})$, is completely determined by the conditional mean, $\mu_{\mathcal{Y}|\mathcal{B}}$, and, perhaps, the common scale parameter, σ .
- ▶ Components of \mathcal{Y} are *conditionally independent*, given \mathcal{B} . Thus, the conditional distribution, $(\mathcal{Y}|\mathcal{B} = \mathbf{b})$, is determined by the (scalar) distribution of each component. Furthermore, the inverse link, \mathbf{g}^{-1} , is determined by a scalar function g^{-1} applied componentwise.

The unscaled conditional density of $(\mathcal{B}|\mathcal{Y} = \mathbf{y})$

- ▶ We observe \mathbf{y} . To make inferences about \mathcal{B} we want the conditional distribution $(\mathcal{B}|\mathcal{Y} = \mathbf{y})$, not $(\mathcal{Y}|\mathcal{B} = \mathbf{b})$.
- ▶ Although the response random variable, \mathcal{Y} , may be discrete or continuous, the random effects, \mathcal{B} , are always a continuous vector-valued random variable.
- ▶ Given values of θ , β and, if used, σ we can evaluate the density of $(\mathcal{B}|\mathcal{Y} = \mathbf{y})$ for any value of \mathbf{b} , but only up to a scale factor.
- ▶ The inverse of the scale factor,

$$\int_{\mathbb{R}^q} [(\mathcal{Y}|\mathcal{B} = \mathbf{b})(\mathbf{y})][\mathcal{B}(\mathbf{b})] d\mathbf{b},$$

is exactly the *likelihood*, $L(\theta, \beta, \sigma^2|\mathbf{y})$, of the parameters given the data.

Penalized (whatever) least squares methods

- ▶ The reason that the PLS problem for determining the conditional modes is relatively easy is because the standard least squares-based methods for fixed-effects models are easily adapted.
- ▶ For linear mixed-models the PLS problem is solved directly. In fact, for LMMs it is possible to determine the conditional modes of the random effects and the conditional estimates of the fixed effects simultaneously.
- ▶ Parameter estimates for generalized linear models (GLMs) are (very efficiently) determined by iteratively re-weighted least squares (IRLS) so the conditional modes in a GLMM are determined by penalized iteratively re-weighted least squares (PIRLS).
- ▶ Nonlinear least squares, used for fixed-effects nonlinear regression, is adapted as penalized nonlinear least squares (PNLS) or penalized iteratively reweighted nonlinear least squares (PIRNLS) for generalized nonlinear mixed models.

The conditional mode of \mathcal{B}

- ▶ To evaluate the likelihood, $L(\theta, \beta, \sigma^2|\mathbf{y})$, at a particular set of parameter values, we first evaluate the *conditional mode* of the random effects,

$$\tilde{\mathbf{b}}(\theta, \beta) = \arg \max_{\mathbf{b}} [(\mathcal{Y}|\mathcal{B} = \mathbf{b})(\mathbf{y})][\mathcal{B}(\mathbf{b})]$$

(σ does not affect the conditional mode.)

- ▶ This optimization problem is easy (well, relatively easy) because it can be expressed as a penalized least squares (PLS) problem. Even better, the PLS problem expressed in terms of *orthogonal random effects*, $\mathcal{U} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I}_q)$, where $\mathcal{B} = \mathbf{T}(\theta)\mathbf{S}(\theta)\mathcal{U}$, has a simple form

$$\tilde{\mathbf{u}}(\theta, \beta) = \arg \min_{\mathbf{u}} \left\| \begin{bmatrix} \mathbf{W}^{1/2} (\mathbf{y} - \mu_{\mathcal{Y}|\mathcal{U}}(\mathbf{u})) \\ \mathbf{u} \end{bmatrix} \right\|^2$$

where the diagonal matrix of weights, \mathbf{W} , depends only on $\mu_{\mathcal{Y}|\mathcal{U}}$ (and, hence, only on η).

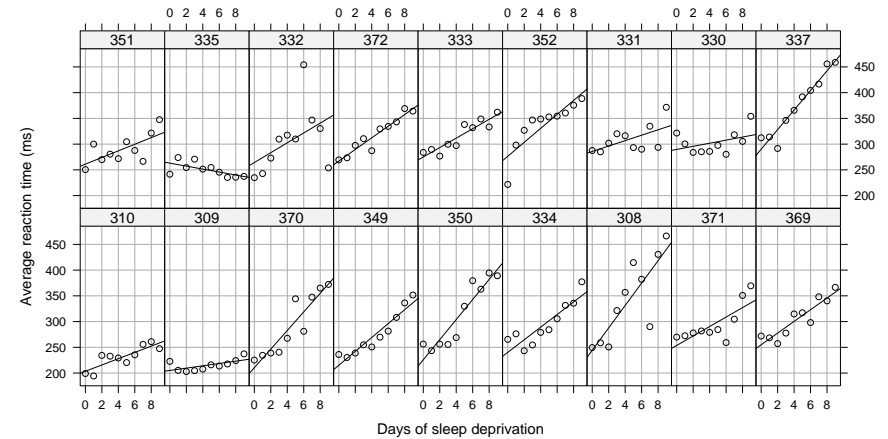
Simple longitudinal data

- ▶ *Repeated measures* data consist of measurements of a response (and, perhaps, some covariates) on several *experimental* (or observational) *units*.
- ▶ Frequently the experimental (observational) unit is *Subject* and we will refer to these units as “subjects”. However, the methods described here are not restricted to data on human subjects.
- ▶ *Longitudinal* data are repeated measures data in which the observations are taken over time.
- ▶ We wish to characterize the response over time within subjects and the variation in the time trends between subjects.
- ▶ Frequently we are not as interested in comparing the particular subjects in the study as much as we are interested in modeling the variability in the population from which the subjects were chosen.

Sleep deprivation data

- ▶ This laboratory experiment measured the effect of sleep deprivation on cognitive performance.
- ▶ There were 18 subjects, chosen from the population of interest (long-distance truck drivers), in the 10 day trial. These subjects were restricted to 3 hours sleep per night during the trial.
- ▶ On each day of the trial each subject's reaction time was measured. The reaction time shown here is the average of several measurements.
- ▶ These data are *balanced* in that each subject is measured the same number of times and on the same occasions.

Reaction time versus days by subject



Comments on the sleep data plot

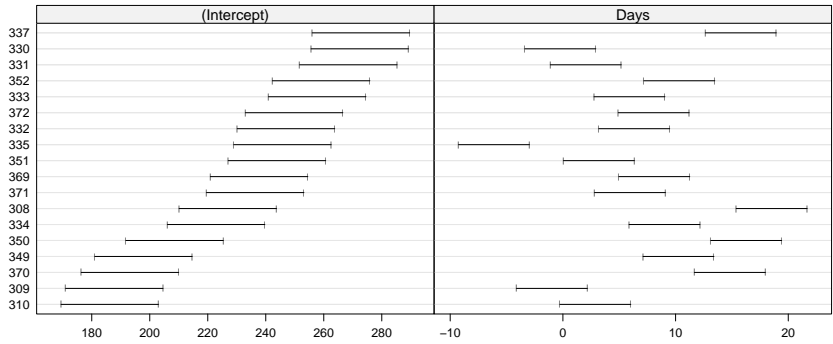
- ▶ The plot is a “trellis” or “lattice” plot where the data for each subject are presented in a separate panel. The axes are consistent across panels so we may compare patterns across subjects.
- ▶ A reference line fit by simple linear regression to the panel's data has been added to each panel.
- ▶ The aspect ratio of the panels has been adjusted so that a typical reference line lies about 45° on the page. We have the greatest sensitivity in checking for differences in slopes when the lines are near $\pm 45^\circ$ on the page.
- ▶ The panels have been ordered not by subject number (which is essentially a random order) but according to increasing intercept for the simple linear regression. If the slopes and the intercepts are highly correlated we should see a pattern across the panels in the slopes.

Assessing the linear fits

- ▶ In most cases a simple linear regression provides an adequate fit to the within-subject data.
- ▶ Patterns for some subjects (e.g. 350, 352 and 371) deviate from linearity but the deviations are neither widespread nor consistent in form.
- ▶ There is considerable variation in the intercept (estimated reaction time without sleep deprivation) across subjects – 200 ms. up to 300 ms. – and in the slope (increase in reaction time per day of sleep deprivation) – 0 ms./day up to 20 ms./day.
- ▶ We can examine this variation further by plotting confidence intervals for these intercepts and slopes. Because we use a pooled variance estimate and have balanced data, the intervals have identical widths.
- ▶ We again order the subjects by increasing intercept so we can check for relationships between slopes and intercepts.

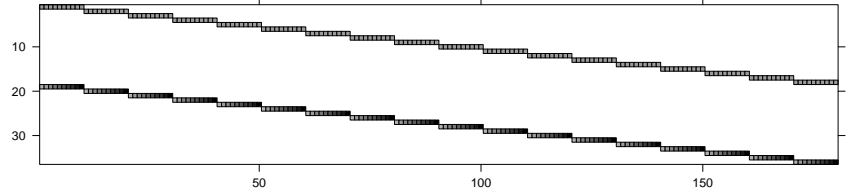
95% conf int on within-subject intercept and slope

A preliminary mixed-effects model



These intervals reinforce our earlier impressions of considerable variability between subjects in both intercept and slope but little evidence of a relationship between intercept and slope.

- ▶ We begin with a linear mixed model in which the fixed effects $[\beta_1, \beta_2]'$ are the representative intercept and slope for the population and the random effects $\mathbf{b}_i = [b_{i1}, b_{i2}]', i = 1, \dots, 18$ are the deviations in intercept and slope associated with subject i .
- ▶ The random effects vector, \mathbf{b} , consists of the 18 intercept effects followed by the 18 slope effects.



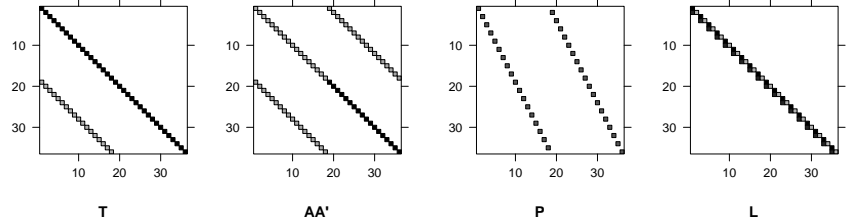
Fitting the model

```
> (fm1 <- lmer(Reaction ~ Days + (Days | Subject),
+   sleepstudy))
```

```
Linear mixed model fit by REML
Formula: Reaction ~ Days + (Days | Subject)
Data: sleepstudy
AIC BIC logLik deviance REMLdev
1756 1775 -871.8 1752 1744
Random effects:
Groups Name Variance Std.Dev. Corr
Subject (Intercept) 612.095 24.7405
Days 35.071 5.9221 0.065
Residual 654.944 25.5919
Number of obs: 180, groups: Subject, 18
Fixed effects:
Estimate Std. Error t value
(Intercept) 251.405 6.825 36.84
Days 10.467 1.546 6.77
Correlation of Fixed Effects:
(Intr)
Days -0.138
```

Terms and matrices

- ▶ The term `Days` in the formula generates a model matrix \mathbf{X} with two columns, the intercept column and the numeric `Days` column. (The intercept is included unless suppressed.)
- ▶ The term `(Days|Subject)` generates a vector-valued random effect (intercept and slope) for each of the 18 levels of the `Subject` factor.



A model with uncorrelated random effects

- ▶ The data plots gave little indication of a systematic relationship between a subject's random effect for slope and his/her random effect for the intercept. Also, the estimated correlation is quite small.
- ▶ We should consider a model with uncorrelated random effects. To express this we use two random-effects terms with the same grouping factor and different left-hand sides. In the formula for an `lmer` model, distinct random effects terms are modeled as being independent. Thus we specify the model with two distinct random effects terms, each of which has `Subject` as the grouping factor. The model matrix for one term is intercept only (1) and for the other term is the column for `Days` only, which can be written `0+Days`. (The expression `Days` generates a column for `Days` and an intercept. To suppress the intercept we add `0+` to the expression; `-1` also works.)

Comparing the models

- ▶ Model `fm1` contains model `fm2` in the sense that if the parameter values for model `fm1` were constrained so as to force the correlation, and hence the covariance, to be zero and the model re-fit we would get model `fm2`.
- ▶ The value 0 to which the correlation is constrained is not on the boundary of the allowable parameter values.
- ▶ In these circumstances a likelihood ratio test and a reference distribution of a χ^2 on 1 degree of freedom is suitable.

```
> anova(fm2, fm1)
Data: sleepstudy
Models:
fm2: Reaction ~ Days + (1 | Subject) + (0 + Days | Subject)
fm1: Reaction ~ Days + (Days | Subject)
   Df    AIC    BIC logLik  Chisq Chi Df Pr(>Chisq)
fm2  5 1762.05 1778.01 -876.02
fm1  6 1763.99 1783.14 -875.99 0.0609     1     0.805
```

A mixed-effects model with independent random effects

```
Linear mixed model fit by REML
Formula: Reaction ~ Days + (1 | Subject) + (0 + Days | Subject)
Data: sleepstudy
   AIC  BIC logLik deviance REMLdev
1754 1770 -871.8   1752   1744
Random effects:
Groups   Name             Variance Std.Dev.
Subject (Intercept) 627.577  25.0515
Subject Days         35.852   5.9876
Residual                653.594  25.5655
Number of obs: 180, groups: Subject, 18
Fixed effects:
              Estimate Std. Error t value
(Intercept)  251.405     6.885   36.51
Days         10.467     1.559    6.71
Correlation of Fixed Effects:
      (Intr)
Days -0.184
```

Conclusions from the likelihood ratio test

- ▶ Because the large p-value indicates that we would not reject `fm2` in favor of `fm1` we prefer the more parsimonious `fm2`.
- ▶ This conclusion is consistent with the AIC (Akaike's Information Criterion) and the BIC (Bayesian Information Criterion) values for which "smaller is better".
- ▶ When evaluating other parameters we will use model `fm2` and an MCMC sample `ss2` from this fitted model (variance and covariance parameters on the transformed scale). The density plots and QQ plots from `ss2` are similar to those from `ss1a` and we do not repeat them here.

Likelihood ratio tests on variance components

- ▶ As for the case of a covariance, we can fit the model with and without the variance component and compare the quality of the fits.
- ▶ The likelihood ratio is a reasonable test statistic for the comparison but the “asymptotic” reference distribution of a χ^2 does not apply because the parameter value being tested is on the boundary.
- ▶ The p-value computed using the χ^2 reference distribution should be conservative (i.e. greater than the p-value that would be obtained through simulation).

```
> fm3 <- lmer(Reaction ~ Days + (1 | Subject), sleepstudy)
> anova(fm3, fm2)
```

```
Data: sleepstudy
```

```
Models:
```

```
fm3: Reaction ~ Days + (1 | Subject)
```

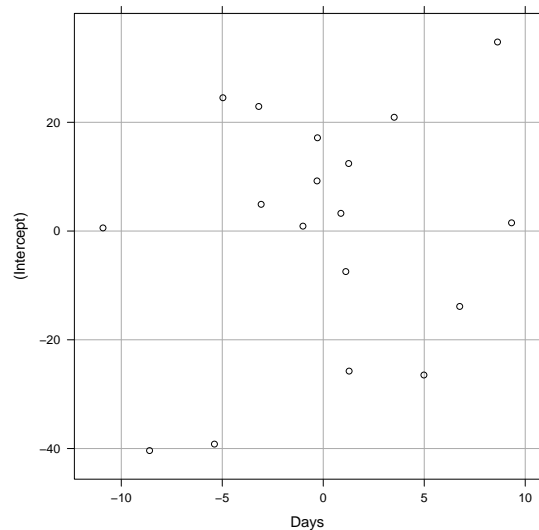
```
fm2: Reaction ~ Days + (1 | Subject) + (0 + Days | Subject)
```

```
      Df    AIC    BIC  logLik  Chisq Chi Df Pr(>Chisq)
```

```
fm3   4 1802.10 1814.87 -897.05
```

```
fm2   5 1762.05 1778.01 -876.02 42.053      1 8.885e-11
```

Scatterplot of the conditional modes



Values of the conditional modes

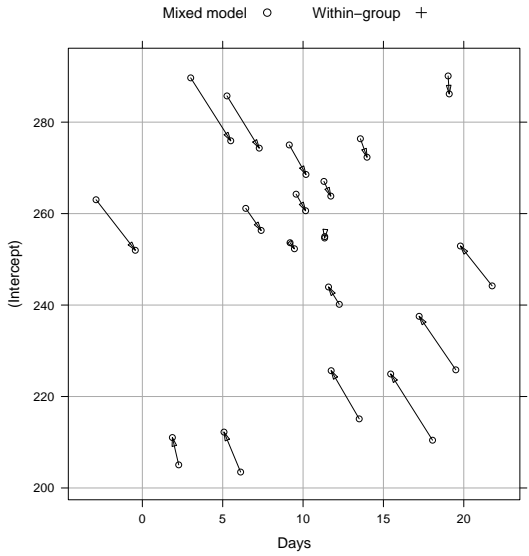
```
> (rr2 <- ranef(fm2))
```

```
$Subject
  (Intercept)    Days
308  1.5138200  9.3232135
309 -40.3749105 -8.5989183
310 -39.1816682 -5.3876346
330  24.5182907 -4.9684965
331  22.9140346 -3.1938382
332   9.2219311 -0.3084836
333  17.1560765 -0.2871973
334  -7.4515945  1.1159563
335   0.5774094 -10.9056435
337  34.7689482  8.6273639
349 -25.7541541  1.2806475
350 -13.8642120  6.7561993
351   4.9156063 -3.0750415
352  20.9294539  3.5121076
369   3.2587507  0.8730251
370 -26.4752098  4.9836365
371   0.9055257 -1.0052631
372  12.4219020  1.2583667
```

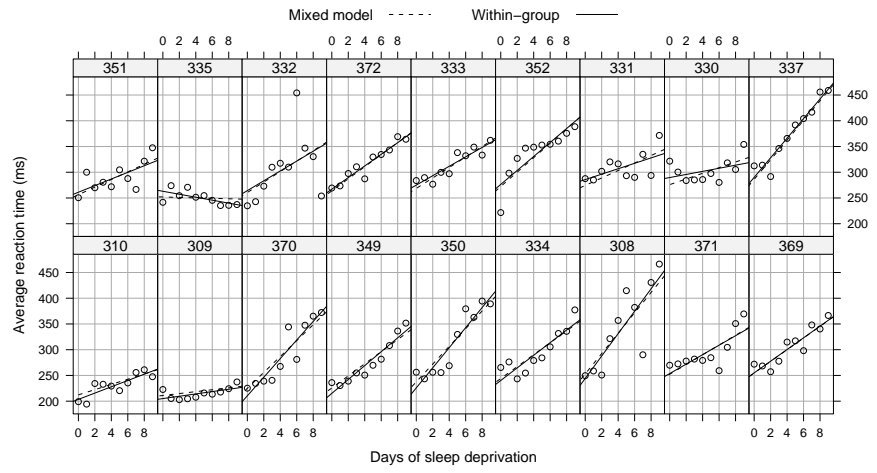
Comparing within-subject coefficients

- ▶ For this model we can combine the conditional modes of the random effects and the estimates of the fixed effects to get conditional modes of the within-subject coefficients.
- ▶ These conditional modes will be “shrunk” towards the fixed-effects estimates relative to the estimated coefficients from each subject’s data. John Tukey called this “borrowing strength” between subjects.
- ▶ Plotting the shrinkage of the within-subject coefficients shows that some of the coefficients are considerably shrunk toward the fixed-effects estimates.
- ▶ However, comparing the within-group and mixed model fitted lines shows that large changes in coefficients occur in the noisy data. Precisely estimated within-group coefficients are not changed substantially.

Estimated within-group coefficients and BLUPs



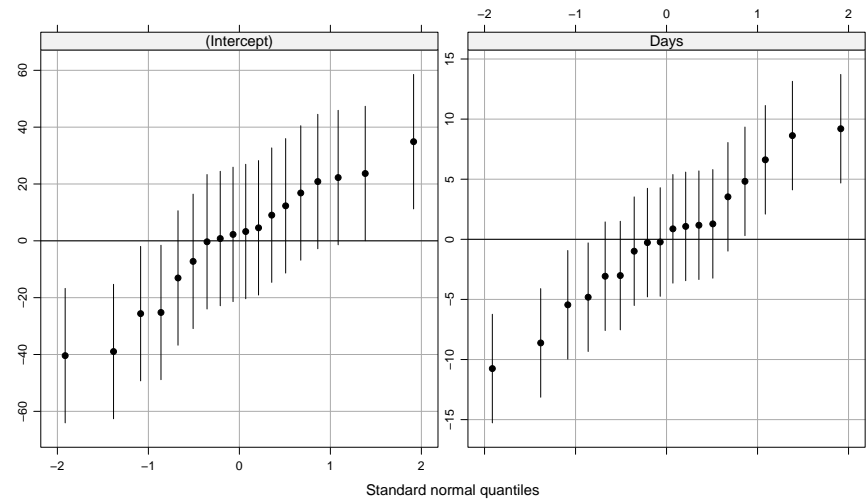
Observed and fitted



Prediction intervals on the random effects

- ▶ For the linear mixed model we can calculate both the means and the variances of the random-effects conditional on the estimated values of the model parameters, which allows us to calculate prediction intervals on the values of individual random effects.
- ▶ We plot the prediction intervals as a normal probability plot so we can see the overall shape of the distribution of the means and which of the random effects are “significantly different” from zero.
- ▶ Note that failure of the conditional means of the random effects to look like a normal (Gaussian) distribution is not terribly alarming. It is the “prior” distribution of the random effects that is assumed to be normal. The conditional means or BLUPs are strongly influenced by the data and may appear non-normal.

Normal probability plot of the prediction intervals



Conclusions from the example

- ▶ Carefully plotting the data is enormously helpful in formulating the model.
- ▶ It is relatively easy to fit and evaluate models to data like these, from a balanced designed experiment.
- ▶ For a linear mixed model the estimates of the fixed effects typically have a symmetric distribution close to a Gaussian distribution.
- ▶ The distribution of the variance components or the covariances are not symmetric, which is why we transform these parameters to a symmetric scale.
- ▶ We use the MCMC sample to create confidence (actually HPD) intervals on the fixed-effects parameters. We could also use the parameter estimates and standard errors.
- ▶ The “estimates” (actually BLUPs) of the random effects can be considered as penalized estimates of these parameters in that they are shrunk towards the origin.
- ▶ Most of the prediction intervals for the random effects overlap zero.

Mean attainment by school

```
> head(patt)
      m attain  n   type
1P 4.425926  54 Primary
2P 5.285714   7 Primary
3P 8.666667   3 Primary
4P 6.285714   7 Primary
5P 4.679245  53 Primary
6P 5.927273  55 Primary

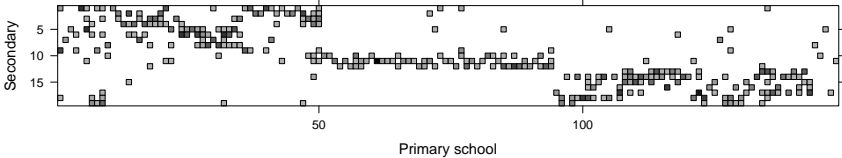
> head(satt)
      m attain  n   type
1S 5.365297 219 Secondary
2S 6.060302 199 Secondary
3S 5.455128 156 Secondary
4S 6.345324 139 Secondary
5S 6.074286 175 Secondary
6S 5.892000 250 Secondary
```

A smaller, non-nested unbalanced example

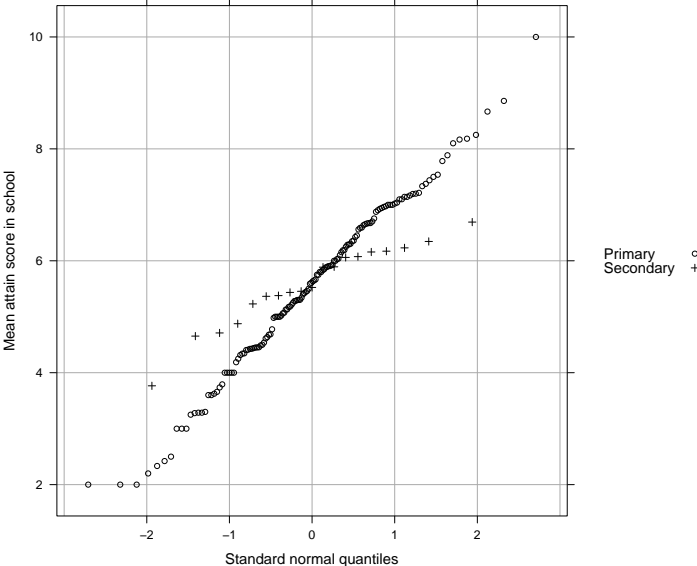
To examine the structure of non-nested, unbalanced observational data we need a smaller example than the 1.3 million observations in the 10 years of grade point scores.

```
> str(ScotsSec)
'data.frame': 3435 obs. of 6 variables:
 $ verbal : num  11 0 -14 -6 -30 -17 -17 -11 -9 -19 ...
 $ attain : num  10 3 2 3 2 2 4 6 4 2 ...
 $ primary: Factor w/ 148 levels "1","2","3","4",...: 1 1 1 1 1 1 1 1 1 1 ...
 $ sex    : Factor w/ 2 levels "M","F": 1 2 1 1 2 2 2 1 1 1 ...
 $ social : num  0 0 0 20 0 0 0 0 0 0 ...
 $ second : Factor w/ 19 levels "1","2","3","4",...: 9 9 9 9 9 9 1 1 9 9

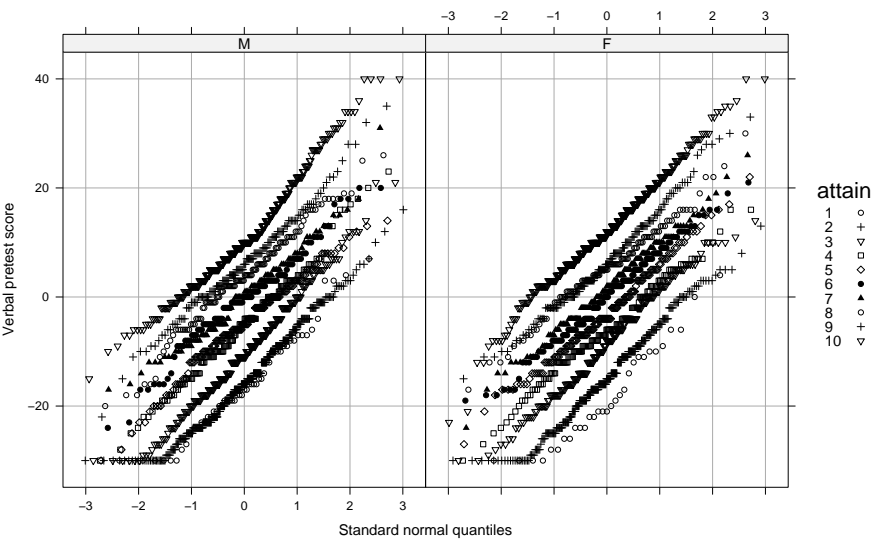
> stab <- xtabs(~second + primary, ScotsSec, sparse = TRUE)
```



Normal probability plot of mean attainment by school



Probability plot of pretest by posttest and sex



An LMM for the secondary school data

```

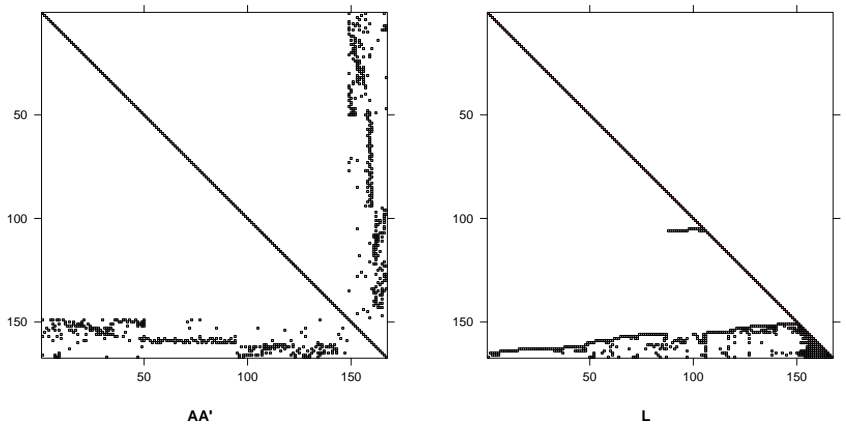
Linear mixed model fit by REML
Formula: attain ~ verbal + sex + (1 | primary) + (verbal | second)
Data: ScotsSec
AIC   BIC logLik deviance REMLdev
14875 14924 -7429   14842   14859
Random effects:
Groups   Name                Variance  Std.Dev.  Corr
primary  (Intercept)             2.7825e-01 0.5274908
second   (Intercept)             2.4862e-02 0.1576754
         verbal              4.0435e-05 0.0063588 0.799
Residual                    4.2434e+00 2.0599441
Number of obs: 3435, groups: primary, 148; second, 19
Fixed effects:
                Estimate Std. Error t value
(Intercept)    5.912760   0.080172   73.75
verbal          0.159249   0.003174   50.17
sexF           0.113791   0.071440    1.59
n = 3435, p = 3, k = 2, n1 = 148, n2 = 19, q1 = 1, q2 = 2,
q = 186
    
```

Reduced LMM for the secondary school data

```

Linear mixed model fit by REML
Formula: attain ~ verbal + sex + (1 | primary) + (1 | second)
Data: ScotsSec
AIC   BIC logLik deviance REMLdev
14872 14909 -7430   14843   14860
Random effects:
Groups   Name                Variance Std.Dev.
primary  (Intercept)             0.276261 0.52561
second   (Intercept)             0.014455 0.12023
Residual                    4.251960 2.06203
Number of obs: 3435, groups: primary, 148; second, 19
Fixed effects:
                Estimate Std. Error t value
(Intercept)    5.919273   0.076151   77.73
verbal          0.159593   0.002778   57.46
sexF           0.115966   0.071463    1.62
n = 3435, p = 3, k = 2, n1 = 148, n2 = 19, q1 = 1, q2 = 1,
q = 167
    
```

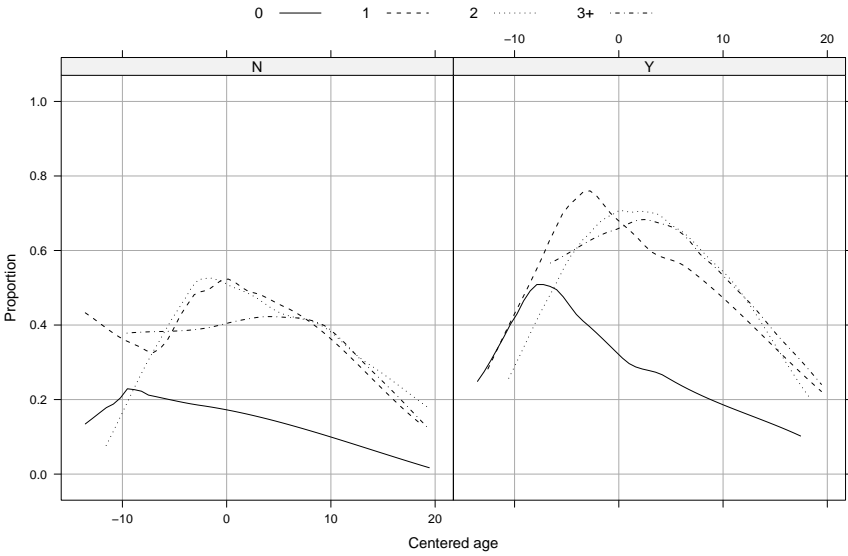
Reordering unbalanced, non-nested data



Generalized linear mixed models

- ▶ A *generalized linear mixed model* (GLMM) is used to model repeated measures data with a non-normal distribution of the response conditional on the value of the linear predictor.
- ▶ Typically the data are binary (0/1) or binomial (k successes out of t trials) or a Poisson count.
- ▶ The formulation of the linear predictor is the same as in the linear mixed model. However, the linear predictor is mapped to the conditional mean via a non-trivial “inverse link” function.
- ▶ The log-likelihood or, equivalently, the deviance of a GLMM does not have an explicit form in this case. As shown later, the marginal likelihood is expressed as an integral that does not have a closed-form solution.
- ▶ We use the Laplace approximation that involves determining the conditional modes of the random effects via PIRLS. A more accurate alternative, adaptive Gauss-Hermite quadrature (AGQ) is being implemented as a Google Summer of Code project. However, AGQ is not universally applicable.

Contraception use versus age by urban and livch



Contraception data

- ▶ One of the data sets in the "mlmRev" package, derived from data files available on the multilevel modelling web site, is from a fertility survey of women in Bangladesh.
- ▶ One of the responses is whether or not the woman currently uses artificial contraception (i.e. a binary response)
- ▶ Covariates included the woman's age (on a centered scale), the number of live children she had, whether she lived in an urban or rural setting, and the district in which she lived.
- ▶ Instead of plotting such data as points, we use the 0/1 response to generate scatterplot smoother curves versus age for the different groups.

Comments on the data plot

- ▶ These observational data are unbalanced (some districts have only 2 observations, some have nearly 120). They are not longitudinal (no “time” variable).
- ▶ Binary responses have low per-observation information content (exactly one bit per observation). Districts with few observations will not contribute strongly to estimates of random effects.
- ▶ Within-district plots will be too imprecise so we only examine the global effects in plots.
- ▶ The comparisons on the multilevel modelling site are for fits of a model that is linear in *age*, which is clearly inappropriate.
- ▶ The form of the curves suggests at least a quadratic in *age*.
- ▶ The urban versus rural differences may be additive.
- ▶ It appears that the *livch* factor could be dichotomized into “0” versus “1 or more”.

Preliminary model fit

```
Generalized linear mixed model fit by the Laplace approximation
Formula: use ~ age + I(age^2) + urban + livch + (1 | district)
Data: Contraception
   AIC   BIC logLik deviance
2389 2433  -1186    2373
Random effects:
  Groups   Name      Variance Std.Dev.
district (Intercept) 0.22586  0.47524
Number of obs: 1934, groups: district, 60
Fixed effects:
             Estimate Std. Error z value Pr(>|z|)
(Intercept) -1.0350725  0.1743606  -5.936 2.91e-09
age          0.0035329  0.0092311   0.383  0.702
I(age^2)     -0.0045623  0.0007252  -6.291 3.15e-10
urbanY       0.6972694  0.1198788   5.816 6.01e-09
livch1       0.8150449  0.1621898   5.025 5.03e-07
livch2       0.9165105  0.1850995   4.951 7.37e-07
livch3+     0.9150209  0.1857689   4.926 8.41e-07
```

Comments on the model fit

- ▶ There is a highly significant quadratic term in `age`.
- ▶ The linear term in `age` is not significant but we retain it because the `age` scale has been centered at an arbitrary (and unknown) value.
- ▶ The `urban` factor is highly significant (as indicated by the plot).
- ▶ Levels of `livch` greater than 0 are significantly different from 0 but may not be different from each other.

Reduced model with dichotomized livch

```
> Contraception$ch <- factor(Contraception$livch != 0,
+   labels = c("N", "Y"))
> print(cm2 <- glmer(use ~ age + I(age^2) + urban + ch +
+   (1 | district), Contraception, binomial), corr = FALSE)
```

```
Generalized linear mixed model fit by the Laplace approximation
Formula: use ~ age + I(age^2) + urban + ch + (1 | district)
Data: Contraception
   AIC   BIC logLik deviance
2385 2419  -1187    2373
Random effects:
  Groups   Name      Variance Std.Dev.
district (Intercept) 0.22470  0.47402
Number of obs: 1934, groups: district, 60
Fixed effects:
             Estimate Std. Error z value Pr(>|z|)
(Intercept) -1.0064262  0.1678949  -5.994 2.04e-09
age          0.0062563  0.0078404   0.798  0.425
I(age^2)     -0.0046354  0.0007163  -6.471 9.73e-11
urbanY       0.6929504  0.1196687   5.791 7.01e-09
chY          0.8603757  0.1473539   5.839 5.26e-09
```

Comparing the model fits

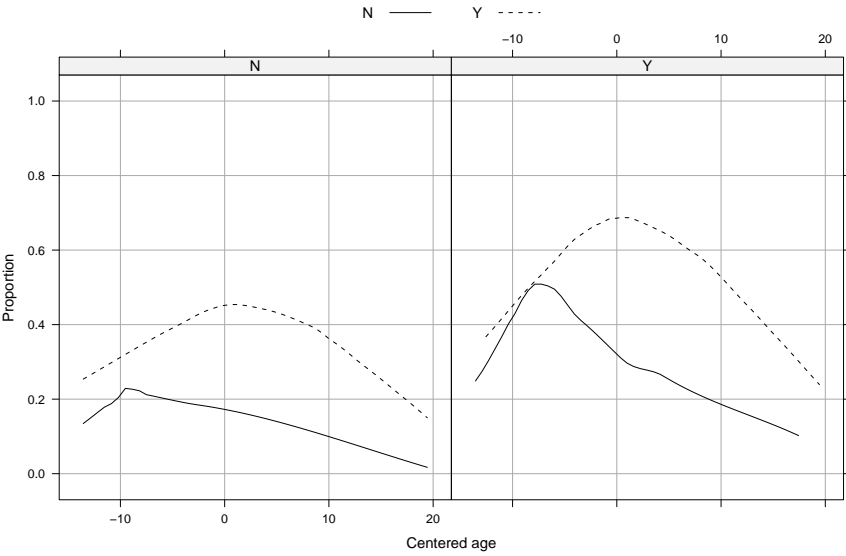
- ▶ A likelihood ratio test can be used to compare these nested models.

```
> anova(cm2, cm1)
```

```
Data: Contraception
Models:
cm2: use ~ age + I(age^2) + urban + ch + (1 | district)
cm1: use ~ age + I(age^2) + urban + livch + (1 | district)
      Df    AIC    BIC logLik  Chisq Chi Df Pr(>Chisq)
cm2   6 2385.2 2418.6 -1186.6
cm1   8 2388.7 2433.3 -1186.4 0.4571    2    0.7957
```

- ▶ The large p-value indicates that we would not reject `cm2` in favor of `cm1` hence we prefer the more parsimonious `cm2`.
- ▶ The plot of the scatterplot smoothers according to live children or none indicates that there may be a difference in the age pattern between these two groups.

Contraception use versus age by urban and ch



Allowing age pattern to vary with ch

Generalized linear mixed model fit by the Laplace approximation

Formula: use ~ age * ch + I(age^2) + urban + (1 | district)

Data: Contraception
 AIC BIC logLik deviance
 2379 2418 -1183 2365

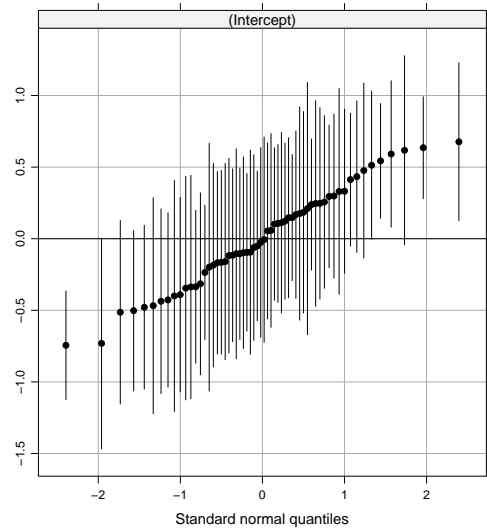
Random effects:
 Groups Name Variance Std.Dev.
 district (Intercept) 0.22306 0.4723
 Number of obs: 1934, groups: district, 60
 Fixed effects:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-1.3233174	0.2144470	-6.171	6.79e-10
age	-0.0472956	0.0218394	-2.166	0.0303
chY	1.2107856	0.2069938	5.849	4.93e-09
I(age^2)	-0.0057572	0.0008358	-6.888	5.64e-12
urbanY	0.7140326	0.1202579	5.938	2.89e-09
age:chY	0.0683522	0.0254347	2.687	0.0072

Correlation of Fixed Effects:

	(Intr) age	chY	I(g^2)	urbanY
age	0.719			
chY	-0.883	-0.793		
I(age^2)	-0.090	0.303	-0.101	

Prediction intervals on the random effects



Extending the random effects

- ▶ We may want to consider allowing a random effect for urban/rural by district. This is complicated by the fact the many districts only have rural women in the study

	district															
urban	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16
N	54	20	0	19	37	58	18	35	20	13	21	23	16	17	14	18
Y	63	0	2	11	2	7	0	2	3	0	0	6	8	101	8	2
	district															
urban	17	18	19	20	21	22	23	24	25	26	27	28	29	30	31	32
N	24	33	22	15	10	20	15	14	49	13	39	45	25	45	27	24

Including a random effect for urban by district

```

Generalized linear mixed model fit by the Laplace approximation
Formula: use ~ age * ch + I(age^2) + urban + (urban | district)
Data: Contraception
AIC BIC logLik deviance
2372 2422 -1177 2354
Random effects:
Groups Name Variance Std.Dev. Corr
district (Intercept) 0.37830 0.61506
urbanY 0.52613 0.72535 -0.793
Number of obs: 1934, groups: district, 60
Fixed effects:
Estimate Std. Error z value Pr(>|z|)
(Intercept) -1.3442631 0.2227667 -6.034 1.60e-09
age -0.0461836 0.0219446 -2.105 0.03533
chY 1.2116527 0.2082372 5.819 5.93e-09
I(age^2) -0.0056514 0.0008431 -6.703 2.04e-11
urbanY 0.7902095 0.1600484 4.937 7.92e-07
age:chY 0.0664682 0.0255674 2.600 0.00933

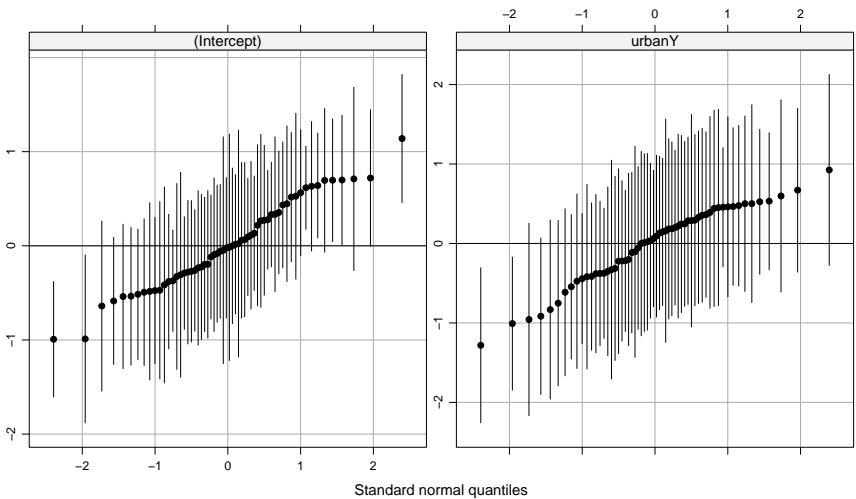
```

```

Correlation of Fixed Effects:
(Intr) age chY I(g^2) urbanY
age 0.696
chY -0.855 -0.792

```

Prediction intervals for the bivariate random effects



Significance of the additional random effect

```

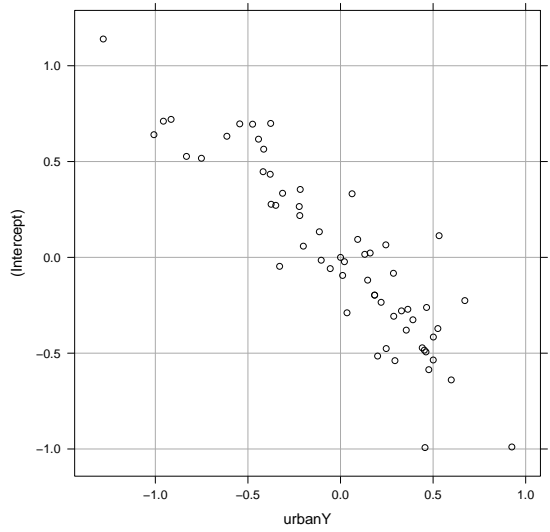
> anova(cm4, cm3)

Data: Contraception
Models:
cm3: use ~ age * ch + I(age^2) + urban + (1 | district)
cm4: use ~ age * ch + I(age^2) + urban + (urban | district)
Df AIC BIC logLik Chisq Chi Df Pr(>Chisq)
cm3 7 2379.2 2418.2 -1182.6
cm4 9 2371.5 2421.6 -1176.8 11.651 2 0.002951

```

- ▶ The additional random effect is highly significant in this test.
- ▶ Most of the prediction intervals still overlap zero.
- ▶ A scatterplot of the random effects shows several random effects vectors falling along a straight line. These are the districts with all rural women or all urban women.

Scatter plot of the conditional modes



Conclusions from the example

- ▶ Again, carefully plotting the data is enormously helpful in formulating the model.
- ▶ Observational data tend to be unbalanced and have many more covariates than data from a designed experiment. Formulating a model is often more difficult than in a designed experiment.
- ▶ A generalized linear model is family, typically `binomial` or `poisson`, is specified as the `family` argument in the call to `glmer`.
- ▶ We use likelihood-ratio tests and z-tests in the model building.

Special case of linear mixed models (cont'd)

- ▶ It is not necessary to solve for $\tilde{\mathbf{u}}(\boldsymbol{\theta})$ and $\hat{\boldsymbol{\beta}}(\boldsymbol{\theta})$. All that is needed for evaluation of the profiled log-likelihood is the penalized residual sum of squares, r^2 , and the determinant

$$|\mathbf{A}\mathbf{A}' + \mathbf{I}| = |\mathbf{L}|^2$$

- ▶ Because \mathbf{L} is triangular, its determinant is simply the product of its diagonal elements.
- ▶ Because $\mathbf{A}\mathbf{A}' + \mathbf{I}$ is positive definite, $|\mathbf{L}|^2 > 0$.
- ▶ The profiled deviance, as a function of $\boldsymbol{\theta}$ only ($\boldsymbol{\beta}$ and σ^2 at their conditional estimates), is

$$d(\boldsymbol{\theta}|\mathbf{y}) = \log(|\mathbf{L}|^2) + n \left(1 + \log(r^2) + \frac{2\pi}{n} \right)$$

Evaluating the likelihood - linear mixed models

- ▶ In the special case of a linear mixed model, where $\boldsymbol{\mu}_{\mathbf{y}|\boldsymbol{\mu}}$ depends linearly on both \mathbf{u} and $\boldsymbol{\beta}$, the conditional mode $\tilde{\mathbf{u}}(\boldsymbol{\theta})$ and the conditional estimate $\hat{\boldsymbol{\beta}}(\boldsymbol{\theta})$ can be determined simultaneously as the solutions to a penalized least squares problem

$$\begin{bmatrix} \tilde{\mathbf{u}}(\boldsymbol{\theta}) \\ \hat{\boldsymbol{\beta}}(\boldsymbol{\theta}) \end{bmatrix} = \arg \min_{\mathbf{u}, \boldsymbol{\beta}} \left\| \begin{bmatrix} \mathbf{y} \\ \mathbf{0} \end{bmatrix} - \begin{bmatrix} \mathbf{A}'\mathbf{P}' & \mathbf{X} \\ \mathbf{I} & \mathbf{0} \end{bmatrix} \begin{bmatrix} \mathbf{u} \\ \boldsymbol{\beta} \end{bmatrix} \right\|^2$$

for which the solution satisfies

$$\begin{bmatrix} \mathbf{P}(\mathbf{A}\mathbf{A}' + \mathbf{I})\mathbf{P}' & \mathbf{P}\mathbf{A}\mathbf{X} \\ \mathbf{X}'\mathbf{A}'\mathbf{P}' & \mathbf{X}'\mathbf{X} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{u}}(\boldsymbol{\theta}) \\ \hat{\boldsymbol{\beta}}(\boldsymbol{\theta}) \end{bmatrix} = \begin{bmatrix} \mathbf{P}\mathbf{A}\mathbf{y} \\ \mathbf{X}'\mathbf{y} \end{bmatrix}$$

- ▶ The Cholesky factor of the system matrix for the PLS problem is

$$\begin{bmatrix} \mathbf{P}(\mathbf{A}\mathbf{A}' + \mathbf{I})\mathbf{P}' & \mathbf{P}\mathbf{A}\mathbf{X} \\ \mathbf{X}'\mathbf{A}'\mathbf{P}' & \mathbf{X}'\mathbf{X} \end{bmatrix} = \begin{bmatrix} \mathbf{L} & \mathbf{0} \\ \mathbf{R}'_{\mathbf{ZX}} & \mathbf{R}'_{\mathbf{X}} \end{bmatrix} \begin{bmatrix} \mathbf{L}' & \mathbf{R}_{\mathbf{ZX}} \\ \mathbf{0} & \mathbf{R}_{\mathbf{X}} \end{bmatrix}$$

- ▶ The dense matrices $\mathbf{R}_{\mathbf{ZX}}$ and $\mathbf{R}_{\mathbf{X}}$ are stored in the `RZX` and `RX` slots, respectively.

REML results

- ▶ Although not often derived in this form, Laird and Ware showed that the REML criterion can be derived as the integral of the likelihood w.r.t. $\boldsymbol{\beta}$.
- ▶ The same techniques as used to evaluate the integral w.r.t. \mathbf{b} can be used to evaluate the integral for the REML criterion. In this case the integral introduces the factor $|\mathbf{R}_{\mathbf{X}}|^2$.
- ▶ The profiled REML deviance, as a function of $\boldsymbol{\theta}$ only (σ at its conditional estimate), is

$$d_R(\boldsymbol{\theta}|\mathbf{y}) = \log(|\mathbf{L}|^2 |\mathbf{R}_{\mathbf{X}}|^2) + (n - p) \left(1 + \log(r^2) + \frac{2\pi}{n - p} \right)$$

Recap

- ▶ For a linear mixed model, even one with a huge number of observations and random effects like the model for the grade point scores, evaluation of the ML or REML profiled deviance, given a value of θ , is straightforward. It involves updating T and S , then updating A , L , R_{ZX} , R_X , calculating the penalized residual sum of squares, r and a couple of determinants of triangular matrices.
- ▶ The profiled deviance can be optimized as a function of θ only. The dimension of θ is usually very small. For the grade point scores there are only three components to θ .

Summary

- ▶ Random variables in an LMM or GLMM
 - \mathcal{Y} the response variable
 - \mathcal{B} the (possibly correlated) random effects
 - \mathcal{U} the orthogonal random effects
- ▶ Parameters in an LMM or GLMM
 - β - fixed-effects coefficients
 - σ - the common scale parameter (not used in some GLMMs)
 - θ - parameters that determine the matrices S and T
- ▶ Matrices used in the definition of an LMM or GLMM
 - X the $n \times p$ model matrix for β
 - Z the $n \times q$ model matrix for u
 - S a $q \times q$ diagonal matrix determined by θ with non-negative diagonal elements
 - T a $q \times q$ unit lower-triangular matrix determined by θ