# Modeling with random effects

Rasmus Waagepetersen
Department of Mathematics
Aalborg University
Denmark

October 2, 2019

## Course topics

- ▶ random effects
- ▶ linear mixed models
- ▶ statistical inference for linear mixed models (including analysis of variance)
- ▶ prediction of random effects
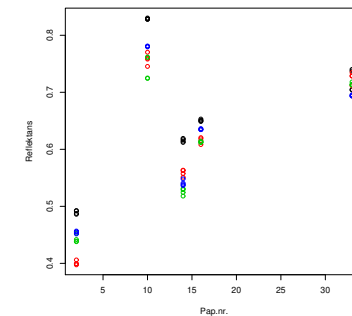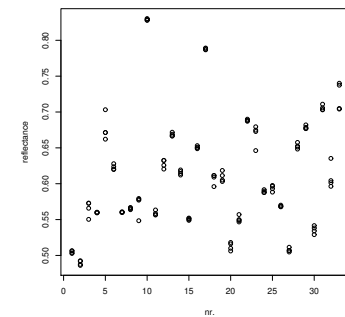- ▶ Implementation in R and SPSS

## Outline

- ▶ examples of data sets
- ▶ random effects models - motivation and interpretation

Next session : details on implementation in R and SPSS

## Reflectance (colour) measurements for samples of cardboard (egg trays) (project at Department of Biotechnology, Chemistry and Environmental Engineering)

Four replications at same position on each cardboard

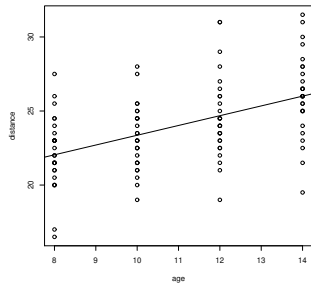For five cardboards: four replications at four positions at each cardboard



Colour variation between/within cardboards ?

## Orthodontic growth curves (repeated measurements/longitudinal data)

Distance (related to jaw size) between pituitary gland and the pterygomaxillary fissure (two distinct points on human skull) for children of age 8-14

Distance versus age:

## Orthodontic growth curves (repeated measurements/longitudinal data)

Distance (related to jaw size) between pituitary gland and the pterygomaxillary fissure (two distinct points on human skull) for children of age 8-14

Distance versus age:

Distance versus age grouped according to child
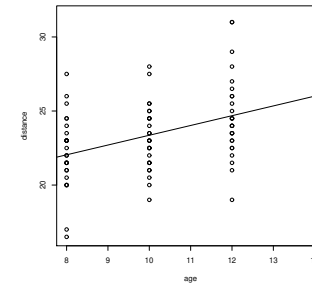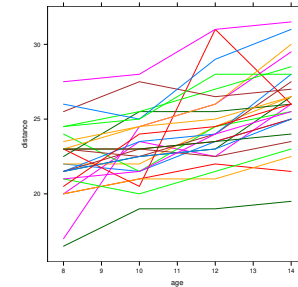



Different intercepts for different children !

Recall: basic aim for statistical analysis of a sample/dataset is to extract information that can be generalized to the population that was sampled.

This perspective in mind when deciding on models for the datasets considered.

## Model for reflectances: one-way anova

Models:

Four replications on each cardboard



$$Y_{ij} = \mu + \epsilon_{ij} \quad i = 1, \ldots, k \ \ j = 1, \ldots, m$$

($k = 34$, $m = 4$) where $\mu$ expectation and $\epsilon_{ij}$ random independent noise

## Model for reflectances: one-way anova

Models:

$$Y_{ij} = \mu + \epsilon_{ij} \quad i = 1, \ldots, k \; j = 1, \ldots, m$$

Four replications on each cardboard

($k = 34$, $m = 4$) where $\mu$ expectation and $\epsilon_{ij}$ random independent noise or

$$Y_{ij} = \mu + \alpha_i + \epsilon_{ij}$$
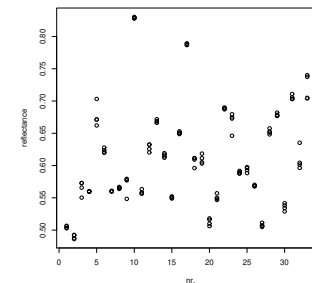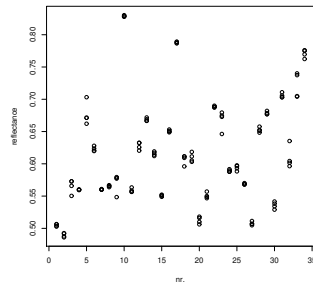
where $\alpha_i$ are fixed unknown parameters

## Model for reflectances: one-way anova

Models:

$$Y_{ij} = \mu + \epsilon_{ij} \quad i = 1, \ldots, k \; j = 1, \ldots, m$$

Four replications on each cardboard

($k = 34$, $m = 4$) where $\mu$ expectation and $\epsilon_{ij}$ random independent noise or

$$Y_{ij} = \mu + \alpha_i + \epsilon_{ij}$$

where $\alpha_i$ are fixed unknown parameters  or

$$Y_{ij} = \mu + U_i + \epsilon_{ij}$$

where $U_i$ are zero-mean random variables independent of each other and of $\epsilon_{ij}$

Which is most relevant ?

## One role of random effects: parsimonious and population relevant models

With fixed effects $\alpha_i$: many parameters ($\mu$, $\sigma^2$, $\alpha_2, \ldots, \alpha_{34}$). Parameters $\alpha_2, \ldots, \alpha_{34}$ not interesting as they just represent intercepts for specific card boards which are individually not of interest.

With random effects: just three parameters ($\mu$, $\sigma^2 = \mathbb{V}\mathrm{ar}\epsilon_{ij}$ and $\tau^2 = \mathbb{V}\mathrm{ar}U_i$).

Hence parsimonious model. Variance parameters interesting for several reasons.

## Second role of random effects: quantify sources of variation

Quantify sources of variation (e.g. quality control): is pulp for paper production too heterogeneous ?

With random effects model

$$Y_{ij} = \mu + U_i + \epsilon_{ij}$$

we have decomposition of variance:

$$\mathbb{V}\mathrm{ar}Y_{ij} = \mathbb{V}\mathrm{ar}U_i + \mathbb{V}\mathrm{ar}\epsilon_{ij} = \tau^2 + \sigma^2$$

Hence we can quantify variation between ($\tau^2$) cardboard pieces and within ($\sigma^2$) cardboard.

Ratio $\gamma = \tau^2/\sigma^2$ is 'signal to noise'.

Proportion of variance

$$\frac{\tau^2}{\sigma^2 + \tau^2} = \frac{\gamma}{\gamma + 1}$$

is called *intra-class correlation*.

High proportion of between cardboard variance leads to high correlation (next slide).

# Third role: modeling of covariance and correlation

Covariances:

$$\mathbb{Cov}[Y_{ij}, Y_{i'j'}] = \begin{cases} 0 & i \neq i' \\ \mathbb{Var}\, U_i & i = i', j \neq j' \\ \mathbb{Var}\, U_i + \mathbb{Var}\, \epsilon_{ij} & i = i', j = j' \end{cases}$$

Correlations:

$$\mathbb{Corr}[Y_{ij}, Y_{i'j'}] = \begin{cases} 0 & i \neq i' \\ \tau^2/(\sigma^2 + \tau^2) & i = i', j \neq j' \\ 1 & i = i', j = j' \end{cases}$$

That is, observations for same cardboard are correlated !

Correct modeling of correlation is important for correct evaluation of uncertainty.

# Fourth role: correct evalution of uncertainty

Suppose we wish to estimate $\mu = \mathbb{E}Y_{ij}$. Due to correlation, observations on same cardboard to some extent redundant.

Estimate is empirical average $\hat{\mu} = \bar{Y}_{..}$. Evaluation of $\mathbb{Var}\bar{Y}_{..}$:

Model erroneously ignoring variation between cardboards

$$Y_{ij} = \mu + \epsilon_{ij}$$

$$\mathbb{Var}\epsilon_{ij} = \sigma^2_{\text{total}} \left[= \sigma^2 + \tau^2\right]$$

Naive variance expression is

$$\mathbb{Var}\bar{Y}_{..} = \frac{\sigma^2_{\text{total}}}{n} \left[= \frac{\sigma^2 + \tau^2}{mk}\right]$$

Correct model with random cardboard effects

$$Y_{ij} = \mu + U_i + \epsilon_{ij},$$

$$\mathbb{Var}\, U_i = \tau^2, \quad \mathbb{Var}\epsilon_{ij} = \sigma^2$$

Correct variance expression is

$$\mathbb{Var}\bar{Y}_{..} = \frac{\tau^2}{k} + \frac{\sigma^2}{mk}$$

With first model, variance is underestimated !

For $\mathbb{Var}\bar{Y}_{..} \to 0$ is it enough that $mk \to \infty$ ?

# Classical balanced one-way ANOVA (analysis of variance)

Decomposition of empirical variance/sums of squares ($i = 1, \ldots, k$, $j = 1, \ldots, m$):

$$SST = \sum_{ij}(Y_{ij} - \bar{Y}_{..})^2 = \sum_{ij}(Y_{ij} - \bar{Y}_{i.})^2 + m\sum_i(\bar{Y}_{i.} - \bar{Y}_{..})^2 = SSE + SSB$$

Expected sums of squares:

$$\mathbb{E}SSE = k(m-1)\sigma^2$$

$$\mathbb{E}SSB = m(k-1)\tau^2 + (k-1)\sigma^2$$

Moment-based estimates:

$$\hat{\sigma}^2 = \frac{SSE}{k(m-1)} \quad \hat{\tau}^2 = \frac{SSB/(k-1) - \hat{\sigma}^2}{m}$$

More complicated formulae in the unbalanced case.

## Hypothesis tests

Fixed effects: $H_0$: $\alpha_1 = \alpha_2 = \cdots = \alpha_k$

$$F = \frac{SSB/(k-1)}{SSE/(k(m-1))}$$

Random effects: $H_0$: $\tau^2 = 0$ Same test-statistic

$$F = \frac{SSB/(k-1)}{SSE/(k(m-1))}$$

Idea: if $\tau^2 = 0$ then $\mathbb{E}SSB/(k-1) = \mathbb{E}SSE/(k(m-1))$.

## Classical implementation in R

For cardboard/reflectance data, $k = 34$ and $m = 4$. `anova()` procedure produces table of sums of squares.

```
> anova(lm(Reflektans~factor(Pap.nr.)))
Analysis of Variance Table

Response: Reflektans
                Df  Sum Sq Mean Sq F value
factor(Pap.nr) 33  0.9009 0.0273   470.7    #SSB
Residuals      102 0.0059 0.00006          #SSE
---
```

Hence $\hat{\sigma}^2 = 0.00006$, $\hat{\tau}^2 = (0.0273 - 0.00006)/4 = 0.00681$.

Biggest part of variation is between cardboard.

## Orthodontic data: classical multiple linear regression in R

```
#fit model with sex specific intercepts and slopes
> ort1=lm(distance~age+age:factor(Sex)+factor(Sex))
> summary(ort1)
Coefficients:
                     Estimate Std. Error t value Pr(>|t|)
(Intercept)           16.3406     1.4162  11.538  < 2e-16 ***
age                    0.7844     0.1262   6.217 1.07e-08 ***
factor(Sex)Female      1.0321     2.2188   0.465    0.643
age:factor(Sex)Female -0.3048     0.1977  -1.542    0.126
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.257 on 104 degrees of freedom
Multiple R-squared: 0.4227,Adjusted R-squared: 0.4061
F-statistic: 25.39 on 3 and 104 DF,  p-value: 2.108e-12
```

Sex and age:Sex not significant !

## Multiple linear regression continued - without interaction

```
> ort2=lm(distance~age+factor(Sex))

Coefficients:
                  Estimate Std. Error t value Pr(>|t|)
(Intercept)       17.70671    1.11221  15.920  < 2e-16 ***
age                0.66019    0.09776   6.753 8.25e-10 ***
factor(Sex)Female -2.32102    0.44489  -5.217 9.20e-07 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '

Residual standard error: 2.272 on 105 degrees of freedom
Multiple R-squared: 0.4095,Adjusted R-squared: 0.3983
F-statistic: 36.41 on 2 and 105 DF,  p-value: 9.726e-13
```
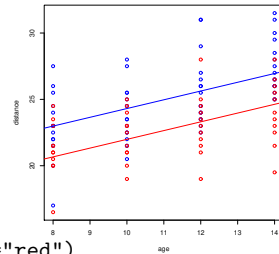
both age and sex significant
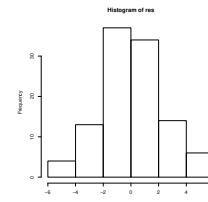
## Multiple linear regression in R III

```
#plot data and two regression lines
col=rep("blue",length(Sex))
col[Sex=="Female"]="red"
plot(distance~age,col=col)
abline(parm[1:2],col="blue")
abline(c(parm[1]+parm[3],parm[2]),col="red")
```
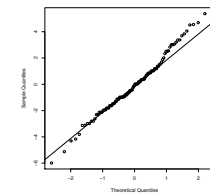
## Multiple linear regression in R IV
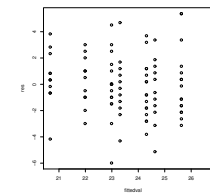
```
res=residuals(ort2)
```

```
hist(res)                qqnorm(res)              fittedval=fitted(ort
                         qqline(res)              plot(res~fittedval)
```
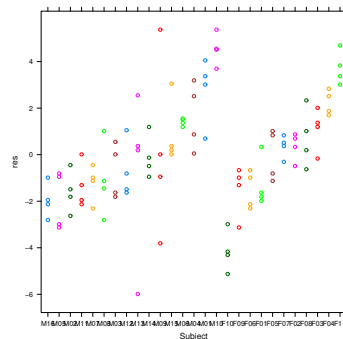
## Multiple linear regression in R V

```
> library(lattice)
> xyplot(res~Subject,groups=Subject)
```



Oups - residuals not independent and identically distributed !
Hence computed $F$-tests not valid.

Problem: subject specific intercepts (and possibly subject specific slopes too)

## Model with subject specific intercepts

```
> ortss=lm(distance~Subject+age+age:factor(Sex)+factor(Sex))
> summary(ortss)

Coefficients: (1 not defined because of singularities)
                    Estimate Std. Error t value Pr(>|t|)
(Intercept)          16.7611     0.6697  25.028  < 2e-16 ***
Subject.L             6.8893     2.9857   2.307  0.02365 *
Subject.Q             0.1675     0.9825   0.170  0.86507
Subject.C             2.7670     1.1527   2.400  0.01873 *
Subject^4             2.8589     0.9497   3.010  0.00350 **
Subject^5            -0.2532     0.7896  -0.321  0.74930
Subject^6            -1.7999     0.8988  -2.003  0.04865 *
Subject^7             0.4857     0.6986   0.695  0.48893
Subject^8             2.4339     0.8380   2.904  0.00477 **
...
Subject^20           -1.3058     0.7276  -1.795  0.07653 .
Subject^21            0.3881     0.6934   0.560  0.57725
Subject^22            2.0115     0.7296   2.757  0.00724 **
Subject^23            1.7772     0.7366   2.413  0.01816 *
Subject^24           -0.7753     0.7025  -1.104  0.27306
Subject^25            1.4231     0.7133   1.995  0.04948 *
Subject^26           -2.1068     0.7292  -2.889  0.00498 **
age                   0.7844     0.0775  10.121 6.44e-16 ***
factor(Sex)Female          NA         NA      NA       NA
age:factor(Sex)Female -0.3048     0.1214  -2.511  0.01410 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 1.386 on 79 degrees of freedom
Multiple R-squared:  0.8345,Adjusted R-squared:  0.7759
F-statistic: 14.23 on 28 and 79 DF,  p-value: < 2.2e-16
```

For each subject an estimate of deviation between the subject's intercept and the first subject's intercept.

In total 27 (!) subject specific estimates.

Each estimate pretty poor (only 4 observations for each subject).

Can not estimate female effect !

Model with subject specific effects may be more correct but is it useful ?

## Mixed model for growth data

$$Y_{ij} = \alpha + \delta_{\text{sex}(i)} + \beta x_{ij} + a_i + b_i x_{ij} + \epsilon_{ij}, \quad i: \text{child}, j: \text{time}$$

Models for coefficients:

- ▶ If interest lies in mean intercept and slope $(\alpha, \beta)$ and sex difference $\delta_s$ but not individual subjects then wasteful to include subject specific fixed effects $a_i$ and $b_i$ (want parsimonious models).

- ▶ Using random effects $a_i$ and $b_i$ with variances $\tau_a^2$ and $\tau_b^2$ allows quantification of population heterogeneity. And only unknown parameters $\alpha$, $\beta$, $\delta_s$, $\tau_a^2$, $\tau_b^2$ and $\sigma^2$ (do not need to estimate $a_i$ and $b_i$)

Back to first role of random effects: parsimonious and meaningful modeling of heterogeneous data. Mixed model: both systematic and random effects.

## Marginal and conditional means of observations

Suppose $a_i \sim N(0, \tau_a^2)$ and $b_i \sim N(0, \tau_b^2)$

Unconditional (marginal) mean of observation:

$$\mathbb{E}[Y_{ij}] = \alpha + \delta_{\text{sex}(i)} + \beta \text{age}_{ij}$$

- i.e. one regression line for each sex (population mean of subject specific lines).

Conditional on $a_i$ and $b_i$:

$$\mathbb{E}[Y_{ij}|a_i, b_i] = [\alpha + a_i] + \delta_{\text{sex}(i)} + [\beta + b_i]\text{age}_{ij}$$

i.e. subject specific lines vary randomly around population mean.

## Mixed model analysis of orthodont data

```
> ort4=lmer(distance~age+Sex+(1|Subject))
> summary(ort4)
Random effects:
 Groups    Name         Variance Std.Dev.
 Subject  (Intercept) 3.2668    1.8074
 Residual               2.0495    1.4316
Number of obs: 108, groups: Subject, 27

Fixed effects:
            Estimate Std. Error t value
(Intercept) 17.70671    0.83391  21.233
age          0.66019    0.06161  10.716
SexFemale   -2.32102    0.76139  -3.048
```

Between subject variance: 3.27, Noise variance: 2.05.

Both age and Sex significant according to Wald-tests (approximate normality of $t$-values).

## Comparison of variances

Total variance: 3.27+2.05=5.32

Similar to estimated residual variance for multiple linear regression model: $5.26 = 2.272^2$.

## Looking at interaction in mixed model framework

```
Formula: distance ~ age * Sex + (1 | Subject)

Random effects:
 Groups    Name         Variance Std.Dev.
 Subject   (Intercept)  3.299    1.816
 Residual               1.922    1.386
Number of obs: 108, groups:  Subject, 27

Fixed effects:
              Estimate Std. Error t value
(Intercept)    16.3406     0.9813  16.652
age             0.7844     0.0775  10.121
SexFemale       1.0321     1.5374   0.671
age:SexFemale  -0.3048     0.1214  -2.511
```

Now interaction significant ($p$=0.012) assuming $t$-value approximately normal.

What is interpretation of interaction ? Does it make sense ?

Note: corresponding model without random effects has much inflated residual variance $5.09 = 2.257^2$ vs. 1.922 for mixed model.

Interaction 'drowns' in large random noise.

## Summary - role of random effects

Models with random effects (mixed models) are useful for:

▶ quantifying different sources of variation

▶ appropriate modeling of variance structure and correlation

▶ correct evalution of uncertainty of parameter estimates

▶ estimation of population variation instead of subject specific characteristics

▶ more parsimonious models (one variance parameter vs. many subject specific fixed effects parameters)

## Exercises

1. Show results regarding covariances and correlations (slide 14) for the $Y_{ij}$ in one-way ANOVA (i.e. the model on slide 12).
2. Analyze the pulp data (brightness of paper pulp in groups given by different operators; from the faraway package) using a one-way anova with random operator effects. Estimate variance components and the intra-class correlation (you may also use output on next slide).

One-way anova for pulp data (4 operators, 5 observations for each operator):

```
> anova(lm(bright~operator,data=pulp))
Analysis of Variance Table

Response: bright
          Df Sum Sq Mean Sq F value  Pr(>F)
operator   3   1.34 0.44667  4.2039 0.02261 * #SSB
Residuals 16   1.70 0.10625                   #SSE
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 '
```

## More exercises

3. compute variances, covariances and correlations of observations from the linear model with random intercepts:

$$Y_{ij} = \alpha + a_i + \beta x_{ij} + \epsilon_{ij}$$

where $\epsilon_{ij} \sim N(0, \sigma^2)$ and $a_i \sim N(0, \tau_a^2)$ and the $\epsilon_{ij}$ and $a_i$ are independent.

4. Continuation of previous exercise. Consider the model fitted on slide 28. What is the proportion of variance due to the error (residual) term ?

5. 5.1 Compute variances, covariances and correlations of observations from the linear model with random slopes:

$$Y_{ij} = \alpha + [\beta + b_i]x_{ij} + \epsilon_{ij}$$

where $\epsilon_{ij} \sim N(0, \sigma^2)$ and $b_i \sim N(0, \tau_b^2)$ and the $\epsilon_{ij}$ and $b_i$ are independent.

5.2 Consider output on next slide. What is the proportion of variance for an observation $Y_{ij}$ explained by the random slopes for different values 8, 10, 12, and 14 of age ?

```
> ort5=lmer(distance~age+Sex+(-1+age|Subject))
> summary(ort5)

Random effects:
 Groups    Name Variance Std.Dev.
 Subject   age  0.026374 0.1624
 Residual       2.080401 1.4424
Number of obs: 108, groups: Subject, 27

Fixed effects:
            Estimate Std. Error t value
(Intercept) 17.43042    0.75066  23.220
age          0.66019    0.06949   9.500
SexFemale   -1.64286    0.68579  -2.396
```

6. compute $\mathbb{V}\mathrm{ar}\,\bar{Y}_{..}$ for one way ANOVA (slide 15).

7. compute the expectations of *SSB* and *SSE* in one-way ANOVA (without loss of generality you may assume $\mu = 0$ since $\mu$ cancels out in the sums of squares).

8. (Design of experiment - one-way ANOVA) Suppose $Y_{ij}$ is outcome of $j$th experiment in $i$th lab, $\tau^2 = 1$ variance between labs and $\sigma^2 = 3$ measurement variance.

   8.1 Suppose we want to make in total 100 experiments. What is then the optimal number of labs that makes $\mathbb{V}\mathrm{ar}\,\bar{Y}_{..}$ minimal?

   8.2 Suppose instead we have available 5000 kr., there is an initial cost of 200 kr. for each lab and subsequently 10 kr. for each experiment. What is then the optimal number $k$ of labs that gives the smallest $\mathbb{V}\mathrm{ar}\,\bar{Y}_{..}$ ?