

# Variance component estimation

Winter School 2022

Jian Zeng

# Outline

- 1 Reminder Linear models
- 2 Balanced designs; one-way models; the sire model
- 3 Balanced designs; one-way models; ANOVA
  - ANOVA: ANalysis Of VAriance
  - Properties of estimators
- 4 Maximum likelihood approaches - ML and REML
  - Why can't we focus only on ANOVA approaches?
  - Maximum likelihood
  - Residual maximum likelihood (REML)
- 5 ML vs REML
- 6 Link to Heritability

# Outline

- 1 Reminder Linear models
- 2 Balanced designs; one-way models; the sire model
- 3 Balanced designs; one-way models; ANOVA
  - ANOVA: ANalysis Of VAriance
  - Properties of estimators
- 4 Maximum likelihood approaches - ML and REML
  - Why can't we focus only on ANOVA approaches?
  - Maximum likelihood
  - Residual maximum likelihood (REML)
- 5 ML vs REML
- 6 Link to Heritability

## Fixed models

Assume we have  $N$  groups of  $T$  individuals each,

$$y_{it} = \mu + \beta_i + e_{it}, \quad i = 1, \dots, N \text{ and } t = 1, \dots, T$$

where

- $\mu$  is an unknown parameter to estimate (fixed effect)
- $\beta_i$  is an unknown parameter to estimate that is constant for all  $t$  at  $i$  fixed (fixed effects)
- $e_{it}$  residuals, with mean  $E(\mathbf{e}) = 0$  and variance-covariance  $\text{Var}(\mathbf{e}) = \sigma_e^2 \mathbf{I}$ .  
 $\text{Var}(y_{it}) = \text{Var}(e_{it}) = \sigma_e^2$

Suppose we have 2 groups ( $i = 1, 2$ ) of 3 individual ( $t = 1, 2, 3$ ).

$$y_{it} = \mu + \beta_i + e_{it}, \text{ or in matrix form } \mathbf{y} = \mu + \mathbf{X}\boldsymbol{\beta} + \mathbf{e},$$

where  $\boldsymbol{\beta}$  is a vector of fixed effects. The levels represents all the levels of interest,

$$\mathbf{y} = \begin{pmatrix} y_{11} \\ y_{12} \\ y_{13} \\ y_{21} \\ y_{22} \\ y_{23} \end{pmatrix}, \boldsymbol{\mu} = \begin{pmatrix} \mu \\ \mu \\ \mu \\ \mu \\ \mu \\ \mu \end{pmatrix}, \mathbf{X} = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{pmatrix}, \boldsymbol{\beta} = \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix}, \mathbf{e} = \begin{pmatrix} e_{11} \\ e_{12} \\ e_{13} \\ e_{21} \\ e_{22} \\ e_{23} \end{pmatrix}$$

## Random models

Assume we have  $N$  groups of  $T$  individuals each,

$$y_{it} = \mu + u_i + e_{it}, \quad i = 1, \dots, N \text{ and } t = 1, \dots, T$$

where

- $\mu$  is an unknown parameter to estimate (fixed effect)
- $(u_1, \dots, u_N)$  is a vector of random values with mean  $E(\mathbf{u}) = 0$  and variance-covariance  $\text{Var}(\mathbf{u}) = \sigma_u^2 \mathbf{I}$  (random effect)
- $e_{it}$  residuals, with mean  $E(\mathbf{e}) = 0$  and variance-covariance  $\text{Var}(\mathbf{e}) = \sigma_e^2 \mathbf{I}$ .  
 $\text{Var}(y_{it}) = \text{Var}(u_i) + \text{Var}(e_{it}) = \sigma_u^2 + \sigma_e^2$

Suppose we have 2 groups ( $i = 1, 2$ ) of 3 individual ( $t = 1, 2, 3$ ).

$$y_{it} = \mu + u_i + e_{it} \text{ or in a matrix form } \mathbf{y} = \mu + \mathbf{Z}\mathbf{u} + \mathbf{e},$$

where  $\mathbf{u}$  is random, the levels are considered drawn from an infinite population of levels,

$$\mathbf{y} = \begin{pmatrix} y_{11} \\ y_{12} \\ y_{13} \\ y_{21} \\ y_{22} \\ y_{23} \end{pmatrix}, \quad \boldsymbol{\mu} = \begin{pmatrix} \mu \\ \mu \\ \mu \\ \mu \\ \mu \\ \mu \end{pmatrix}, \quad \mathbf{Z} = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{pmatrix}, \quad \mathbf{u} = \begin{pmatrix} u_1 \\ u_2 \end{pmatrix}, \quad \mathbf{e} = \begin{pmatrix} e_{11} \\ e_{12} \\ e_{13} \\ e_{21} \\ e_{22} \\ e_{23} \end{pmatrix}$$

## Mixed models

Mixed models (MM) contain both fixed and random factors

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e}$$

where

- $\mathbf{y}$  vector of observed dependent values, with mean  $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$
- $\boldsymbol{\beta}$  vector of unknown parameters to estimate (fixed effects)
- $\mathbf{u}$  vector of unknown random effects, with mean  $E(\mathbf{u}) = 0$  and variance-covariance  $\text{Var}(\mathbf{u}) = \mathbf{G}$  (usually  $\mathbf{G} = \sigma_u^2 \mathbf{I}$ )
- $\mathbf{e}$  vector of residuals, with mean  $E(\mathbf{e}) = 0$  and variance-covariance  $\text{Var}(\mathbf{e}) = \mathbf{R}$  (usually  $\mathbf{R} = \sigma_e^2 \mathbf{I}$ ),
- $\mathbf{X}$  and  $\mathbf{Z}$  are design matrices

$$\text{Var}(\mathbf{y}) = \mathbf{V} = \mathbf{Z}\mathbf{G}\mathbf{Z}^\top + \mathbf{R}$$

## Example

Fitting a mean, unrelated sires, uncorrelated errors

$$\mathbf{y} = \boldsymbol{\mu} + \mathbf{Z}\mathbf{u} + \mathbf{e}$$

$$\mathbf{V} = \text{Var}(\mathbf{y}) = \mathbf{Z}\mathbf{Z}^\top \sigma_u^2 + \mathbf{I}\sigma_e^2$$

$$\text{If } \mathbf{Z} = \begin{pmatrix} 1 & 0 \\ 1 & 0 \\ 1 & 0 \\ 0 & 1 \\ 0 & 1 \\ 0 & 1 \end{pmatrix}, \text{ then } \mathbf{Z}\mathbf{Z}^\top = \begin{pmatrix} 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 1 & 1 & 1 \end{pmatrix}$$

$$\text{And we get } \mathbf{V} = \begin{pmatrix} \sigma_u^2 + \sigma_e^2 & \sigma_u^2 & \sigma_u^2 & 0 & 0 & 0 \\ \sigma_u^2 & \sigma_u^2 + \sigma_e^2 & \sigma_u^2 & 0 & 0 & 0 \\ \sigma_u^2 & \sigma_u^2 & \sigma_u^2 + \sigma_e^2 & 0 & 0 & 0 \\ 0 & 0 & 0 & \sigma_u^2 + \sigma_e^2 & \sigma_u^2 & \sigma_u^2 \\ 0 & 0 & 0 & \sigma_u^2 & \sigma_u^2 + \sigma_e^2 & \sigma_u^2 \\ 0 & 0 & 0 & \sigma_u^2 & \sigma_u^2 & \sigma_u^2 + \sigma_e^2 \end{pmatrix}$$

## Aim

## GBLUP reminder

$$\mathbf{y} = \mathbf{1}_n\mu + \mathbf{Z}\mathbf{g} + \mathbf{e}$$

breeding values  $\mathbf{g} \sim \mathcal{N}(0, \mathbf{G}\sigma_a^2)$ , where  $\mathbf{G}$  is the genomic relationship matrix.  
GBLUP solves the following system of equations

$$\begin{bmatrix} \mathbf{1}'_n\mathbf{1}_n & \mathbf{1}'_n\mathbf{Z} \\ \mathbf{Z}'\mathbf{1}_n & \mathbf{Z}'\mathbf{Z} + \mathbf{G}^{-1}\frac{\sigma_e^2}{\sigma_a^2} \end{bmatrix} \begin{bmatrix} \hat{\mu} \\ \hat{\mathbf{g}} \end{bmatrix} = \begin{bmatrix} \mathbf{1}'_n\mathbf{y} \\ \mathbf{Z}'\mathbf{y} \end{bmatrix}$$

BLUP: estimate the mean  $\mu$  and predict breeding values  $\mathbf{g}$ , all based on known variance

Aim of these lectures: estimate variance components (estimate  $\lambda$ )

Two ways to estimate variance components:

- ANOVA: ANalysis Of VAriance
- Maximum Likelihood approaches



# Outline

- 1 Reminder Linear models
- 2 **Balanced designs; one-way models; the sire model**
- 3 Balanced designs; one-way models; ANOVA
  - ANOVA: ANalysis Of VAriance
  - Properties of estimators
- 4 Maximum likelihood approaches - ML and REML
  - Why can't we focus only on ANOVA approaches?
  - Maximum likelihood
  - Residual maximum likelihood (REML)
- 5 ML vs REML
- 6 Link to Heritability

## Sire model

Consider a sire model with  $s$  unrelated sires and  $n$  progeny per sire; one-way model (one random effect: sire), with a balanced design (same number of observations per sire).

$$y_{it} = \mu + u_i + e_{it}$$

$$\text{Var}(y_{it}) = \sigma_u^2 + \sigma_e^2$$

$$\text{Cov}(y_{it}, y_{ik}) = \sigma_u^2$$

$$i = 1, \dots, s$$

$$t = 1, \dots, n$$

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e}$$

$$\text{Var}(\mathbf{y}) = \mathbf{Z}\mathbf{Z}^\top \sigma_u^2 + \mathbf{I}\sigma_e^2$$

$$\text{Var}(\mathbf{y}_i) = \mathbf{Z}_i\mathbf{Z}_i^\top \sigma_u^2$$

$\mathbf{X}$ : matrix of one column of  $sn$  ones

$$\boldsymbol{\beta} = \mu$$

```
## Simulating data
> set.seed(123)
> mu = 1.2; s = 10; n = 20
> sigmau = 1; sigmae = 0.2
> ui = rnorm(s, sd=sigmau)
> eit = rnorm(n*s, sd=sigmae)
> y = matrix(NA_real_, nrow=n*s, ncol=1)
> for(i in 1:s){
+   ind = (n*(i-1)+1) : (n*(i-1)+n)
+   y[ind] = mu + ui[i] + eit[ind]
+ }

## create a grouping factor
> grp = factor(rep(1:s, each = n))

## create Z
> Z = matrix(0, nrow=s*n, ncol=s)
> for(i in 1:s){
+   ind = which(grp == i)
+   Z[indic, i] = 1
+ }
```

# Outline

- 1 Reminder Linear models
- 2 Balanced designs; one-way models; the sire model
- 3 Balanced designs; one-way models; ANOVA**
  - ANOVA: ANalysis Of VAriance
  - Properties of estimators
- 4 Maximum likelihood approaches - ML and REML
  - Why can't we focus only on ANOVA approaches?
  - Maximum likelihood
  - Residual maximum likelihood (REML)
- 5 ML vs REML
- 6 Link to Heritability

# Outline

- 1 Reminder Linear models
- 2 Balanced designs; one-way models; the sire model
- 3 Balanced designs; one-way models; ANOVA**
  - ANOVA: ANalysis Of VAriance
  - Properties of estimators
- 4 Maximum likelihood approaches - ML and REML
  - Why can't we focus only on ANOVA approaches?
  - Maximum likelihood
  - Residual maximum likelihood (REML)
- 5 ML vs REML
- 6 Link to Heritability

## ANOVA table - general framework

| Source of variation      | df                | SS    | MS         | E(MS)         |
|--------------------------|-------------------|-------|------------|---------------|
| Mean                     | $df_M = 1$        | $SSM$ | $SSM/df_M$ | $E(SSM/df_M)$ |
| Between sires            | $df_B = s - 1$    | $SSB$ | $SSB/df_S$ | $E(SSB/df_S)$ |
| Within sires (residuals) | $df_W = s(n - 1)$ | $SSW$ | $SSW/df_W$ | $E(SSW/df_W)$ |
| Total                    | $N$               | $SST$ |            |               |

df: degrees of freedom

SS: sum of squares

MS: Mean square (mean of SS)

E(MS): Expectation of MS

## Sums of squares

We decompose  $SST$  into a mean, a between and a within family component:

$$\begin{aligned}
 SST &= \sum_i \sum_t y_{it}^2 = \sum_i \sum_t (y_{it} - \bar{y} + \bar{y})^2 \\
 &= \sum_i \sum_t (y_{it} - \bar{y})^2 + N\bar{y}^2 \\
 &= \sum_i \sum_t [(y_{it} - \bar{y}_i) + (\bar{y}_i - \bar{y})]^2 + N\bar{y}^2 \\
 &= \sum_i \sum_t [(y_{it} - \bar{y}_i)^2 + 2(y_{it} - \bar{y}_i)(\bar{y}_i - \bar{y}) + (\bar{y}_i - \bar{y})^2] + N\bar{y}^2
 \end{aligned}$$

$\bar{y}$  is the grand mean,  $\bar{y}_i$  is the family mean.

The middle term is equal to zero by definition of a mean:  $\sum_t (y_{it} - \bar{y}_i) = 0$ .

The last term is independent of  $j$  so  $\sum_i \sum_t (\bar{y}_i - \bar{y})^2 = n \sum_i (\bar{y}_i - \bar{y})^2$

### Decomposition of SS

$$\begin{aligned}
 SST &= N\bar{y}^2 + n \sum_i (\bar{y}_i - \bar{y})^2 + \sum_i \sum_t (y_{it} - \bar{y}_i)^2 \\
 &= SS_M + SS_B + SS_W
 \end{aligned}$$

## ANOVA table - our sire model

Using

$$SSM = \left( \sum_i \sum_t y_{it} \right)^2 / (sn)$$

$$SSB = \sum_i \left( \sum_t y_{it} \right)^2 / n - \left( \sum_i \sum_t y_{it} \right)^2 / (sn)$$

$$SSW = \sum_i \sum_t y_{it}^2 - SSB - SSM$$

$$SST = \sum_i \sum_t y_{it}^2$$

we get

| Source of variation      | df         | SS    | MS                   | E(MS)                               |
|--------------------------|------------|-------|----------------------|-------------------------------------|
| Mean                     | 1          | $SSM$ | $SSM/1$              | $N\mu^2 + n\sigma_u^2 + \sigma_e^2$ |
| Between sires            | $s - 1$    | $SSB$ | $B = SSB/(s - 1)$    | $n\sigma_u^2 + \sigma_e^2$          |
| Within sires (residuals) | $s(n - 1)$ | $SSW$ | $W = SSW/(s(n - 1))$ | $\sigma_e^2$                        |
| Total                    | $N = sn$   | $SST$ |                      |                                     |

cf lecture notes for calculations and proofs

## Sum of Squares - Sire model

$$SST = \sum_i \sum_t y_{it}^2$$

$$SSM = (\sum_i \sum_t y_{it})^2 / (sn)$$

$$SSB = \sum_i (\sum_t y_{it})^2 / n - (\sum_i \sum_t y_{it})^2 / (sn)$$

$$SSW = \sum_i \sum_t y_{it}^2 - SSB - SSM$$

```

> SST = sum(y^2); SST
[1] 422.0067

> SSM = sum(y)^2 / (s*n); SSM
[1] 279.1494

> SSB = 0
> for(i in 1:s){
+   ind = (n*(i-1)+1) : (n*(i-1)+n)
+   yit = sum(y[ind])
+   SSB = SSB + yit^2/n
+ }
> SSB = SSB - SSM; SSB
[1] 133.8501

> SSW = sum(y^2) - SSB - SSM; SSW
[1] 9.007134

```



## Estimation of $\sigma_u^2$ and $\sigma_e^2$

### Principle of ANOVA

"Equate SS of analysis of variance to their expected values, giving a set of equations that are linear in the variance components to be estimated"

For the one-way design: two equations, two unknowns:

$$\begin{aligned} B &= n\sigma_u^2 + \sigma_e^2 \\ W &= \sigma_e^2 \end{aligned}$$

Hence,

$$\begin{aligned} \widehat{\sigma_e^2} &= W \\ \widehat{\sigma_u^2} &= (B - W)/n \end{aligned}$$

```
> B = SSB/(s-1); B
[1] 14.87224

> W = SSW/(s*(n-1)); W
[1] 0.04740597

> c(sqrt(W), sigmae)
[1] 0.2177291 0.2000000

> c(sqrt((B - W)/n), sigmau)
[1] 0.8609538 1.0000000
```

# Outline

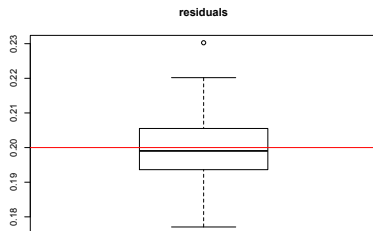
- 1 Reminder Linear models
- 2 Balanced designs; one-way models; the sire model
- 3 Balanced designs; one-way models; ANOVA**
  - ANOVA: ANalysis Of VAriance
  - **Properties of estimators**
- 4 Maximum likelihood approaches - ML and REML
  - Why can't we focus only on ANOVA approaches?
  - Maximum likelihood
  - Residual maximum likelihood (REML)
- 5 ML vs REML
- 6 Link to Heritability

## 50 repeats - no(?) Bias

```

> u = e = NULL
> for(iter in 1:50){
+   ui = rnorm(s,sd=sigmau)
+   eit = rnorm(n*s, sd=sigmae)
+
+   y = mu + Z %* % ui + eit
+
+   SST = sum(y^2)
+   SSM = t(y) %*% X %*%
+   solve(t(X) %*% X) %*%
+   t(X) %*% y
+   SSB = t(y) %*% Z %*%
+   solve(t(Z)%*%Z) %*%
+   t(Z) %*% y - SSM
+   SSW = sum(y^2) - SSB - SSM
+
+   B = SSB/(s-1); B
+   W = SSW/(s*(n-1)); W
+
+   e = c(e, sqrt(W))
+   u = c(u, sqrt((B - W)/n))
+ }

```



## Properties of estimators

### 1 Unbiased

$$E(\widehat{\sigma_e^2}) = E(W) = \sigma_e^2$$

$$\begin{aligned} E(\widehat{\sigma_u^2}) &= E(B - W)/n = E(B)/n - E(W)/n \\ &= (\sigma_u^2 + \sigma_e^2/n) - \sigma_e^2/n \\ &= \sigma_u^2 \end{aligned}$$

- 2 Minimum Variance: Estimates have minimum variance among all possible unbiased estimators. True for normal and non-normal data
- 3 Distribution: Under normality, only the estimate of the residual variance has a  $\chi^2$  distribution

#### 4 Sampling variances

Using that  $SSW \sim \sigma_e^2 \chi^2(df_W)$ , and  $\text{Var}(SSW) = 2\sigma_e^4 df_W$ , we have

$$\begin{aligned}\text{Var}(\widehat{\sigma_e^2}) &= \text{Var}(W) = \text{Var}(SSW)/df_W^2 \\ &= df_W 2\sigma_e^4 / df_W^2 \\ &= 2\sigma_e^4 / df_W\end{aligned}$$

$$\begin{aligned}\text{Var}(\widehat{\sigma_u^2}) &= \text{Var}((B - W)/n) \\ &= [\text{Var}(B) + \text{Var}(W)] / n^2 \\ &= [2E(B)^2 / df_B + 2E(W)^2 / df_W] / n^2 \\ &= [2(n\sigma_u^2 + \sigma_e^2)^2 / df_B + 2\sigma_e^4 / df_W] / n^2 \\ &= (2/n^2) [(n\sigma_u^2 + \sigma_e^2)^2 / (s - 1) + \sigma_e^4 / (s(n - 1))]\end{aligned}$$

$$\text{Var}(\hat{t}) = \text{Var}((B - W)/(B + (n - 1)W))$$

```
> var(e^2)
[1] 1.728204e-05
> 2*sigmae^4/(s*(n-1))
[1] 1.684211e-05
>
> var(u^2)
[1] 0.2296663
> 2/n^2 * ( (n*sigmau^2+sigmae^2)^2/(s-1) + sigmae^4/(s*(n-1)) )
[1] 0.223112
```

# Outline

- 1 Reminder Linear models
- 2 Balanced designs; one-way models; the sire model
- 3 Balanced designs; one-way models; ANOVA
  - ANOVA: ANalysis Of VAriance
  - Properties of estimators
- 4 **Maximum likelihood approaches - ML and REML**
  - Why can't we focus only on ANOVA approaches?
  - Maximum likelihood
  - Residual maximum likelihood (REML)
- 5 ML vs REML
- 6 Link to Heritability

# Outline

- 1 Reminder Linear models
- 2 Balanced designs; one-way models; the sire model
- 3 Balanced designs; one-way models; ANOVA
  - ANOVA: ANalysis Of VAriance
  - Properties of estimators
- 4 Maximum likelihood approaches - ML and REML**
  - Why can't we focus only on ANOVA approaches?**
  - Maximum likelihood
  - Residual maximum likelihood (REML)
- 5 ML vs REML
- 6 Link to Heritability

## Problems with Unbalanced Designs

- SS can be partitioned in many ways
  - ▶ fit effect  $\alpha$  before  $\beta$
  - ▶ fit effect  $\beta$  before  $\alpha$
- no obvious SS (or other quadratic form) to estimate variance components from
- SS (MS) are not *orthogonal* (independent)
- using an ANOVA approach gives biased estimates of variance components for a mixed model



## Methods proposed for unbalanced designs

- Henderson's (1953) methods I, II, and III
  - ▶ Essentially Least Squares methods
  - ▶ Problems with mixed models
- Maximum Likelihood (ML)
  - ▶ Unified procedure for estimating fixed effects and variance components
  - ▶ Desirable asymptotic properties
  - ▶ Bias in variance components
- Residual (restricted) Maximum Likelihood (REML)
  - ▶ Similar to ANOVA for balanced designs
  - ▶ No bias due to loss in degrees of freedom for fitting fixed effects

# Outline

- 1 Reminder Linear models
- 2 Balanced designs; one-way models; the sire model
- 3 Balanced designs; one-way models; ANOVA
  - ANOVA: ANalysis Of VAriance
  - Properties of estimators
- 4 **Maximum likelihood approaches - ML and REML**
  - Why can't we focus only on ANOVA approaches?
  - **Maximum likelihood**
  - Residual maximum likelihood (REML)
- 5 ML vs REML
- 6 Link to Heritability

## Log-likelihood - general case

Our model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e}$  can also be written as a generalised linear model (GLM)

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}, \text{ where } \mathbf{e} \sim \mathcal{N}_n(\mathbf{0}, \mathbf{V})$$

with  $\mathbf{V} = \mathbf{Z}\mathbf{G}\mathbf{Z}^\top + \mathbf{R}$ , where usually  $\mathbf{G} = \sigma_u^2 \mathbf{I}$  and  $\mathbf{R} = \sigma_e^2 \mathbf{I}$ .

The likelihood of such model is

$$L(\boldsymbol{\beta}, \sigma_u^2, \sigma_e^2) = \left(\frac{1}{2\pi}\right)^{N/2} |\mathbf{V}|^{-1/2} \exp\left(\frac{1}{2}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})\right)$$

The log-likelihood in the general case is

$$\ell(\boldsymbol{\beta}, \sigma_u^2, \sigma_e^2) = -\frac{1}{2} \left[ \log(|\mathbf{V}|) + (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \right]$$

The log-likelihood for the sire model is

$$\ell(\boldsymbol{\beta}, \sigma_u^2, \sigma_e^2) = -\frac{1}{2} \left[ \log(|\mathbf{V}|) + (\mathbf{y} - \mathbf{1}\boldsymbol{\mu})^\top \mathbf{V}^{-1}(\mathbf{y} - \mathbf{1}\boldsymbol{\mu}) \right],$$

and  $\mathbf{V} = \mathbf{Z}\mathbf{G}\mathbf{Z}^\top + \mathbf{R} = \mathbf{Z}\mathbf{Z}^\top \sigma_u^2 + \mathbf{I}\sigma_e^2$ .

## Maximum Likelihood approach - one-way model

It follows (Searle, Linear Models, page 418) that,

$$\ell(\beta, \sigma_u^2, \sigma_e^2) = -\frac{1}{2} [s \log(n\sigma_u^2 + \sigma_e^2) + s(n-1) \log(\sigma_e^2) + SSW/\sigma_e^2 + SSB/(n\sigma_u^2 + \sigma_e^2) + sn(\bar{y} - \mu)^2/(n\sigma_u^2 + \sigma_e^2)]$$

### Maximum Likelihood Estimation (MLE)

Taking differentials with respect to  $\mu, \sigma_u^2$  and  $\sigma_e^2$ , we obtain

$$\begin{aligned}\hat{\mu} &= \bar{y} \\ \hat{\sigma}_u^2 &= \left( \left( \frac{s-1}{s} \right) B - W \right) / n \\ \hat{\sigma}_e^2 &= W\end{aligned}$$

with the condition that  $((s-1)/s)B \geq W$

=> classic estimate of the mean, same estimate for  $\sigma_e^2$  as ANOVA, but **biased estimate of  $\sigma_u^2$**  ( $(B-W)/n$  for anova)

```
> c(sum(y)/(n*s), mu)
[1] 1.275278 1.200000

> c(((s-1)/s*B-W)/n, sigmau)
[1] 0.7988936 1.0000000

> c(sqrt(W), sigmae)
[1] 0.1923142 0.2000000
```

# Outline

- 1 Reminder Linear models
- 2 Balanced designs; one-way models; the sire model
- 3 Balanced designs; one-way models; ANOVA
  - ANOVA: ANalysis Of VAriance
  - Properties of estimators
- 4 Maximum likelihood approaches - ML and REML**
  - Why can't we focus only on ANOVA approaches?
  - Maximum likelihood
  - **Residual maximum likelihood (REML)**
- 5 ML vs REML
- 6 Link to Heritability

## ML vs REML

### Maximum Likelihood (ML)

$$\ell = - \left[ (\log(|\mathbf{V}|) + (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top \mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\boldsymbol{\beta})) \right] / 2$$

'Determinant of a variance matrix, plus a weighted sum of squares of residuals'

### Residual (or restricted) Maximum Likelihood (REML)

$$\ell_R = - \left[ (\log(|\mathbf{V}|) - \log(|\mathbf{X}^\top \mathbf{V}^{-1} \mathbf{X}|^{-1})) + (\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}})^\top \mathbf{V}^{-1}(\mathbf{y} - \mathbf{X}\hat{\boldsymbol{\beta}}) \right] / 2$$

$-\log(|\mathbf{X}^\top \mathbf{V}^{-1} \mathbf{X}|^{-1})$ : variance term associated with the estimation of  $\hat{\boldsymbol{\beta}}$   
"penalty term"

Where does that come from?

## Penalty term

### Marginal model

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{e}$$

- $E(\mathbf{y}) = \mathbf{X}\boldsymbol{\beta}$
- $\text{Var}(\mathbf{y}) = \mathbf{V} = \mathbf{Z}\mathbf{G}\mathbf{Z}^\top \sigma_u^2 + \mathbf{R}\sigma_e^2$
- $\mathbf{y} \sim \mathcal{N}(\mathbf{X}\boldsymbol{\beta}, \mathbf{Z}\mathbf{G}\mathbf{Z}^\top \sigma_u^2 + \mathbf{R}\sigma_e^2)$

Linear combinations of  $\mathbf{y}$  have a non-negative variance ( $\mathbf{V}$  is 'non-negative definite')

If we knew the matrix  $\mathbf{V}$ , then

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{V}^{-1} \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{V}^{-1} \mathbf{y} \text{ (weighted least squares)}$$

$$\text{Var}(\hat{\boldsymbol{\beta}}) = (\mathbf{X}^\top \mathbf{V}^{-1} \mathbf{X})^{-1}$$

Side note: the model  $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \mathbf{e}$  is more restrictive than the marginal model: variances of both  $\mathbf{e}$  and  $\mathbf{u}$  are non-negative.

## Estimation/computation

How to maximise the likelihood?

- Maximise likelihood = maximise log-likelihood
- Many methods, e.g.,
  - ▶ Derivative free
  - ▶ Expectation Maximisation
  - ▶ Using second differentials

### Fixed effects and REML

Usual estimates are,

$$\hat{\beta} = (\mathbf{X}^T \hat{\mathbf{V}}^{-1} \mathbf{X})^{-1} \mathbf{X}^T \hat{\mathbf{V}}^{-1} \mathbf{y}$$

These estimates of fixed effects are **not** maximum likelihood estimates!!!

- The likelihood was optimised independent of the fixed effects
- ML properties for the estimates of fixed effects do not apply – > no LRT, e.g. use Wald test



## REML for balanced one-way model

$$\ell_R = -\frac{1}{2} [(s-1) \log(n\sigma_u^2 + \sigma_e^2) + s(n-1) \log(\sigma_e^2) + SSW/\sigma_e^2 + SSB/(n\sigma_u^2 + \sigma_e^2)]$$

We solve the partial derivatives  $\frac{\partial \ell_R}{\partial \sigma_e^2}$  and  $\frac{\partial \ell_R}{\partial \sigma_u^2}$  and obtain

## REML estimates

$$\begin{aligned}\widehat{\sigma_e^2} &= W \\ \widehat{\sigma_u^2} &= (B - W)/n\end{aligned}$$

Same as ANOVA estimates (balanced, one way)!

## Estimates

- If  $B > W$ , then the maximum likelihood estimates are identical to the ANOVA estimates
- If  $B < W$  (negative ANOVA estimates), then

$$\widehat{\sigma_u^2} = 0$$

$$\begin{aligned}\widehat{\sigma_e^2} &= (w_1 W + w_2 B)/(w_1 + w_2) \\ &= (SSB + SSW)/(ns - 1) \\ &= TSS/(N - 1)\end{aligned}$$

where the weights  $w_i = 1/\text{variance}$ .

$$w_1 = 1/(2\sigma_e^2/(s(n-1)))$$

$$w_2 = 1/(2\sigma_e^4/(s-1))$$

# Outline

- 1 Reminder Linear models
- 2 Balanced designs; one-way models; the sire model
- 3 Balanced designs; one-way models; ANOVA
  - ANOVA: ANalysis Of VAriance
  - Properties of estimators
- 4 Maximum likelihood approaches - ML and REML
  - Why can't we focus only on ANOVA approaches?
  - Maximum likelihood
  - Residual maximum likelihood (REML)
- 5 **ML vs REML**
- 6 Link to Heritability

## Balanced designs; on-way models

$$\hat{\mu} = \bar{y}$$

$$\hat{\sigma}_u^2 = \left( \left( \frac{s-1}{s} \right) B - W \right) / n$$

$$\hat{\sigma}_e^2 = W$$

```
> c(sum(y)/(n*s), mu)
[1] 1.275278 1.200000
```

```
> c((s-1)/s*B-W)/n, sigmau)
[1] 0.7988936 1.0000000
```

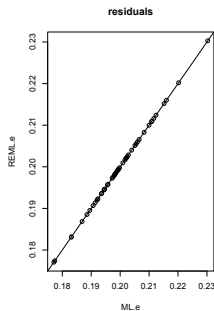
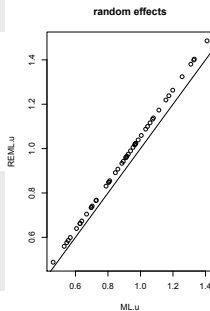
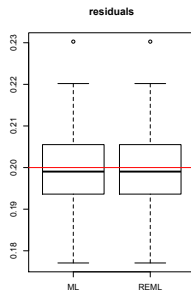
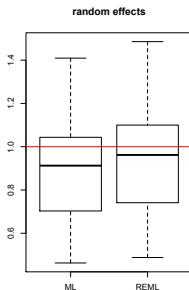
```
> c(sqrt(W), sigmae)
[1] 0.1923142 0.2000000
```

$$\widehat{\sigma}_e^2 = W$$

$$\widehat{\sigma}_u^2 = (B - W) / n$$

```
> c(sqrt(W), sigmae)
[1] 0.2177291 0.2000000
```

```
> c(sqrt((B - W)/n), sigmau)
[1] 0.8609538 1.0000000
```



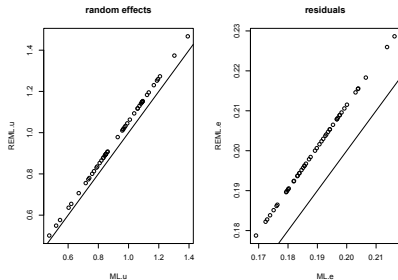
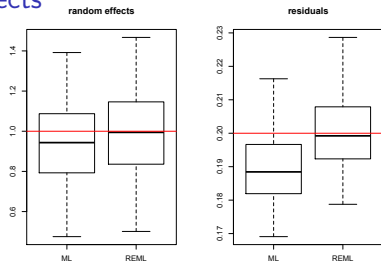
## Balanced designs; on-way models; fixed effects

```

> library(lme4)
> p = 20
> X = matrix(rnorm(s*n*p), nrow=s*n)
> X = scale(X)
> beta = rnorm(p, sd=0.1)

> ML.u=ML.e=
+ REML.u=REML.e= vector("numeric", length=50)
> for(iter in 1:50){
+   ui = rnorm(s,sd=sigmau)
+   eit = rnorm(n*s, sd=sigmae)
+
+   y = mu + X%*% beta + Z %*% ui + eit
+
+   a = lmer(y~ X+ (1|grp), REML=FALSE) # ML
+   b = lmer(y~ X+ (1|grp), REML=TRUE) # REML
+
+   ML.u[iter] = as.data.frame(VarCorr(a))[1,5]
+   ML.e[iter] = as.data.frame(VarCorr(a))[2,5]
+   REML.u[iter] = as.data.frame(VarCorr(b))[1,5]
+   REML.e[iter] = as.data.frame(VarCorr(b))[2,5]
+ }

```



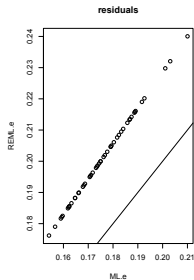
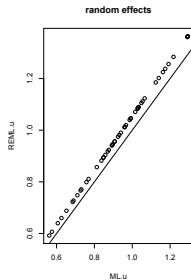
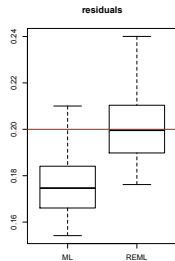
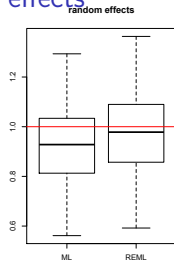
## Unbalanced designs; one-way models; fixed effects

```

> n = sample(5:20,s, replace=TRUE)
> grp=factor(rep(1:s, n))

> ML.u=ML.e=
+ REML.u=REML.e= vector("numeric", length=50)
> for(iter in 1:50){
+   ui = rnorm(s,sd=sigmau)
+   eit = rnorm(n*s, sd=sigmae)
+
+   y = mu + X%*% beta + Z %*% ui + eit
+
+   a = lmer(y~ X+ (1|grp), REML=FALSE) # ML
+   b = lmer(y~ X+ (1|grp), REML=TRUE) # REML
+
+   ML.u[iter] = as.data.frame(VarCorr(a))[1,5]
+   ML.e[iter] = as.data.frame(VarCorr(a))[2,5]
+   REML.u[iter] = as.data.frame(VarCorr(b))[1,5]
+   REML.e[iter] = as.data.frame(VarCorr(b))[2,5]
+ }

```



# Outline

- 1 Reminder Linear models
- 2 Balanced designs; one-way models; the sire model
- 3 Balanced designs; one-way models; ANOVA
  - ANOVA: ANalysis Of VAriance
  - Properties of estimators
- 4 Maximum likelihood approaches - ML and REML
  - Why can't we focus only on ANOVA approaches?
  - Maximum likelihood
  - Residual maximum likelihood (REML)
- 5 ML vs REML
- 6 Link to Heritability

## Genome-Wide Complex Trait Analysis (GCTA)

$$\mathbf{y} = \mathbf{1}_n\mu + \mathbf{g} + \mathbf{e}, \quad (1)$$

$\text{Var}(\mathbf{g}) = \mathbf{G}\sigma_a^2$  with  $\mathbf{G}$  the relatedness (GRM) matrix;  $\mathbf{V} = \text{Var}(\mathbf{y}) = \mathbf{G}\sigma_g^2 + \mathbf{I}\sigma_e^2$ .

$\sigma_g^2$  and  $\sigma_e^2$  estimated by REML;  $h^2 = \frac{\sigma_g^2}{\text{Var}(\mathbf{y})}$ .

*limitation of R:* lmer: random effects are assumed independent; you cannot input the GRM. lmer4qt1 seems to be a recent alternative.

In any case, R is not advised for ML/REML analysis with big datasets (too slow)

We use GCTA