

Chapter 5

Variance Components Estimation

In applied prediction of breeding values using BLUP animal models, variance components for all random effects are required as input. These variance components must be estimated from the data. In more detail, given the assumed linear mixed effect model

$$y = Xb + Za + e \quad (5.1)$$

where y is a vector of length N of observations, b is a vector of length p of fixed effects, a is a vector of length q of random breeding values and e is a vector of length N of random errors. The matrices X and Z are design matrices linking the corresponding effects to the observations. As part of the model definition, the variances of the random effects are defined as

$$\begin{aligned} \text{var}(a) &= A\sigma_a^2 \\ \text{var}(e) &= I\sigma_e^2 \end{aligned} \quad (5.2)$$

In (5.2) σ_a^2 and σ_e^2 are the variance components that must be estimated from the data. The material presented in this chapter is based on [Essl, 1987] and [Searle et al., 1992] and it shows different methods how variance components for different models can be estimated.

5.1 Estimation Of Genetic Components

For each trait that should be considered in an aggregate genotype, the first thing to be analysed is whether the trait has any genetic component. Because only

traits with a detectable genetic component can be used for improving a population on the genetic level. The genetic component quantifies the part that is passed from parents to offspring. Hence from a livestock breeding point of view, the ratio between the genetic variability (quantified by σ_a^2) and the phenotypic variability (measured by σ_p^2) is important and is termed as **heritability** (h^2).

$$h^2 = \frac{\sigma_a^2}{\sigma_p^2} \quad (5.3)$$

One first method that we want to introduce is based on the very well-known statistical technique called **analysis of variance** (ANOVA). ANOVA is shown in the next subsection for a simple application of estimating the repeatability. Later this can be generalized to the estimation of genetic components.

5.2 Estimation Of Repeatability

The term **repeatability** indicates how similar repeated measurements of the same quantity are. For example, if we measure the same trait on any given animal several times, the measurements are expected to vary. But because the measurements are done on the same animal, the variability is probably smaller compared to measurements from different animals. This phenomenon can be quantified by a ratio of variance components which is called repeatability.

The computation of the repeatability is shown using the following example dataset from 10 randomly selected bulls. From each bull the shoulder height is measured three times.

Table 5.1: Repeated Measurements of Shoulder Height in cm

Bull	M1	M2	M3
1	135	136	134
2	129	130	128
3	135	133	136
4	127	127	125
5	126	129	129
6	128	129	128
7	127	132	130
8	129	128	125
9	126	125	127
10	132	131	134

Now we want to check whether the measurements for the same bull have a smaller variability compared to measurements from different bulls. We first

create a plot which might already give us some indications.

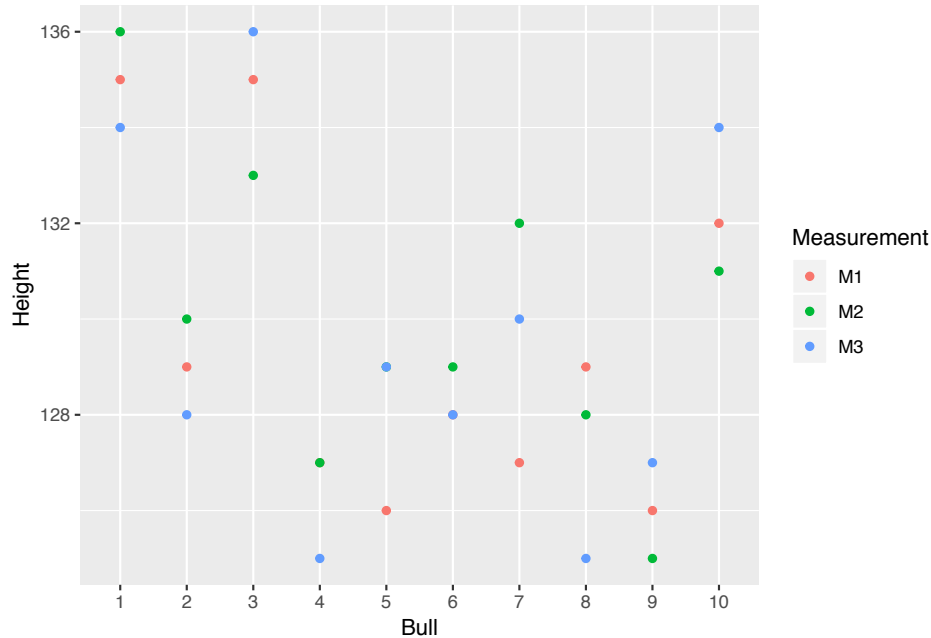


Figure 5.1: Repeated Measurements of Shoulder Height for Ten Bulls

From Figure 5.1 alone, it is difficult to say whether measurements for the same animal are more similar than measurements from different animals. We use the following model to provide a quantitative answer for the previously posed question.

$$y_{ij} = \mu + t_i + \epsilon_{ij} \quad (5.4)$$

where

- y_{ij} measurement j of animal i
- μ expected value of y
- t_i deviation of y_{ij} from μ attributed to animal i
- ϵ_{ij} measurement error

5.2.1 Estimation

Given the definition of t_i and ϵ_{ij} as random effects, the following relationships hold

- $E(t_i) = 0$
- $\sigma_t^2 = E(t_i^2)$: variance component of total variance (σ_y^2) which can be attributed to the t -effects

- $E(\epsilon_{ij}) = 0$
- $\sigma_\epsilon^2 = E(\epsilon_{ij}^2)$: variance component attributed to ϵ -effects
- $\sigma_y^2 = \sigma_t^2 + \sigma_\epsilon^2$

The repeatability w is defined as the following ratio between variance components

$$w = \frac{\sigma_t^2}{\sigma_t^2 + \sigma_\epsilon^2} \quad (5.5)$$

The variance components σ_t^2 and σ_ϵ^2 are estimated using an analysis of variance. The result of such an analysis is shown in the following table.

```
##                Df Sum Sq Mean Sq F value    Pr(>F)
## Bull             9  286.7   31.85   13.85 8.74e-07 ***
## Residuals       20   46.0    2.30
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Based on the theory of analysis of variance the expected values of the mean sum of squares can be equated to the following variance components.

Effect	$E(\text{MeanSq})$
Bull	$\sigma_\epsilon^2 + n * \sigma_t^2$
Error	σ_ϵ^2
Total	$\sigma_\epsilon^2 + \frac{N-n}{N-1} * \sigma_t^2$

where n is the number of measurement per bull and N is the total number of measurements.

The numeric values of the compute **Mean Sq** values are now taken as estimates for the respective variance components. Therefore

$$\hat{\sigma}_\epsilon^2 = 2.3$$

and

$$\hat{\sigma}_t^2 = \frac{31.85 - 2.3}{3} = 9.85$$

The estimated repeatability can now be computed as

$$\hat{w} = \frac{\hat{\sigma}_t^2}{\hat{\sigma}_t^2 + \hat{\sigma}_\epsilon^2} = 0.81$$

5.3 Estimation Of Sire Variance

The technique of estimating variance components using ANOVA can also be applied to a data set where offspring performance records are grouped by their sires using a sire model. From the statistical point of view a sire model is a linear mixed effects model for each observation, the effect of the sire is expressed by a random effect. In matrix vector notation this model can be written as

$$y = Xb + Zs + e \quad (5.6)$$

where y is a vector of length N of observations, b is a vector of length p of fixed effects, s is a vector of length r with random sire effects and e is a vector of length N of random error terms. The matrices X and Z are incidence matrices for b and s , linking the respective effects to the observations. An example of such a data set is used in Problem 1 of Exercise 2.

The variance component σ_s^2 for the random sire component s is estimated the same way as shown in subsection 5.2 using an ANOVA table. For the sire model the ANOVA table has the following structure

Effect	Degrees of Freedom	Sum Sq	Mean Sq	$E(\text{Mean Sq})$
Sire ($s b$)	$r - 1$	$SSQ(s b)$	$SSQ(s b)/(r - 1)$	$\sigma_e^2 + k * \sigma_s^2$
Residual (e)	$N - r$	$SSQ(e)$	$SSQ(e)/(N - r)$	σ_e^2

where

$$SSQ(s|b) = SSQ(sb) - SSQ(b)$$

$$SSQ(sb) = \sum_{i=1}^r \left[\left(\sum_{j=1}^{n_i} y_{ij} \right)^2 / n_i \right]$$

$$SSQ(b) = \left(\sum_{i=1}^r \sum_{j=1}^{n_i} y_{ij} \right)^2 / N$$

$$SSQ(e) = SSQ(y) - SSQ(sb)$$

$$SSQ(y) = \sum_{i=1}^r \sum_{j=1}^{n_i} y_{ij}^2$$

$$k = \frac{1}{r - 1} \left[N - \frac{\sum_{i=1}^r n_i^2}{N} \right]$$

with r the number of sires and n_i the number of progeny for sire i .

The numeric computation of estimating σ_s^2 and σ_e^2 is the topic of Problem 1 of Exercise 2. The dataset that is used in Exercise 2 is a simplified version where only certain genetic relationships occur and where the number of environmental effects are kept at a very low number. To address the higher complexity of real-world datasets obtained in the field, other methods have been developed. Furthermore the ANOVA-based techniques when applied to real data can produce negative estimates for variance components. Because variance components are on a quadratic scale, they cannot be negative and from negative variances, the standard deviations are not defined in the scope of real numbers. Hence negative variance component estimates are outside of the parameter domain.

5.4 Development Of Further Methods

In this subsection, we focus on methods which are still used today. The currently used methods for variance components estimation are either based on Likelihood approaches or are the result of some Bayesian procedure.

5.4.1 Maximum Likelihood

The first maximum likelihood approach to estimate variance components for linear mixed effects models was developed by [Hartley and Rao, 1967]. As the term `maximum likelihood` implies it, the presented method is based on the likelihood L where L is defined as

$$L(\theta) = f(y|\theta) \quad (5.7)$$

where θ is the vector of all unknown parameters to be estimated. For the linear mixed effect model

$$y = Xb + Zu + e$$

and under the assumption of the data being normally distributed, [Hartley and Rao, 1967] specify L as

$$L(\theta) = (2\pi)^{-1/2n} \sigma^{-n} |H|^{-1/2} * \exp \left\{ -\frac{1}{2\sigma^2} (y - Xb)^T H (y - Xb) \right\} \quad (5.8)$$

where $\text{var}(y) = H\sigma^2 = Z^T G Z + R$ with $\text{var}(u) = I\sigma_u^2$ and $\text{var}(e) = R = I\sigma^2$. The maximum likelihoods for σ_u^2 and σ^2 are the values that maximize the function likelihood function L . It has to be noted that in (5.8) not only the

variance components, but also the fixed effects b are unknown. These must also be estimated from the data.

The maximization of L is done by taking the partial derivatives of $\lambda = \log L$ with respect to all unknown parameters. Then these partial derivatives are set to 0 and the resulting solutions are taken as maximum likelihood estimates.

The problem with the just described maximum likelihood approach is that the unknown fixed effects b have to be estimated at the same time. As a consequence of that the maximum likelihood estimates of the variance components depend on b . This is considered as an undesirable property. The solution for this problem was developed by [Patterson and Thompson, 1971] and is called **Restricted Maximum Likelihood** (REML). In REML the observations y are transformed as Sy and Qy with the following properties

- (i) The matrix S has rank $n - t$ and the matrix Q has rank t
- (ii) The result of the two transformations are independent, that means $cov(Sy, Qy) = 0$ which is met when $SHQ^T = 0$
- (iii) The matrix S is chosen such that $E(Sy) = 0$ which means $SX = 0$
- (iv) The matrix QX is of rank t , so that every linear function of the elements of Qy estimate a linear function of b .

From (i) and (ii) it follows that the likelihood L of y is the product of the likelihoods of Sy and Qy that means

$$\lambda = \lambda' + \lambda''$$

Suitable matrices S and Q are given by

$$S = I - X(X^T X)^{-1} X^T$$

and

$$Q = X^T H^{-1}.$$

With these transformations, the variance components σ^2 and σ_u^2 can be estimated by maximizing λ' which is the logarithm of the likelihood of Sy and is independent of any influence of the fixed effects b . Based on this property, REML is the de-facto standard for variance components estimation in applied livestock breeding. The R-package `pedigreemm` can be used to get estimates for variance components using either Maximum Likelihood (ML) or REML.

5.5 Bayesian Procedures

Theoretical foundations for using Bayesian methods in animal breeding were laid by [Gianola and Fernando, 1986]. These foundations spanned more than

just the topic of variance components. A detailed implementation scheme using Gibbs sampling for datasets and models originating in the area of livestock breeding was first described by [Wang et al., 1994]. Parameter estimates were obtained from their respective marginal posterior distribution. These marginal posterior distributions were obtained from applying the Gibbs sampling scheme to the joint posterior distribution. The prior distributions for the variance components are set to be scale-free inverted chi-square distributions.

Because the availability of widely used and tested software implementing Bayesian procedures is limited, these procedures are not used in practical livestock breeding.