

The R-Package 'synbreed'

Valentin Wimmer*

September 2, 2010

Abstract

This document gives an introduction to the R-Package 'synbreed' which contains tools and methods for plant and animal breeding. The goal is the creation of an analysis pipeline for genomic selection. This comprises tools for genotypic, phenotypic and pedigree data. The steps of a typical analysis are presented in this document. This starts with the coding of the marker data, followed by the construction of relationship matrixes according to pedigree or genomic relationship matrixes, i.e. according to vanRaden (2008).

Keywords: synergistic plant and animal breeding, kinship, pedigree, marker data, animal model

1 Introduction

In this document the steps of an analysis pipeline for genotypic and phenotypic data in plant or animal breeding with the R-package 'synbreed' are presented. As an illustration, a simulated data set for maize called `maize` which is part of the package is used. To load `maize` data set, use

```
> library(synbreed)
> data(maize)
```

This data set contains genotypic and phenotypic data for 1250 genotypes of maize. When loading the `maize` data set, in fact the following data sets are loaded in the workspace

maize.geno This is a `data.frame` containing the marker data of 696 biallelic SNP-markers for the 1250 genotypes. The first column of the data set contains the ID for the identification of the genotypes. This variable should be used for the merge with the phenotypic data. The

*Author of correspondence: Institute for plant breeding, Technical University of Munich, Emil-Ramann-Str. 4, 85354 Freising, Germany, Email: Valentin.Wimmer@wzw.tum.de

marker data is coded with 0/1 and because of simulated data no missing values are present. Note that the coding does not contain any information if an allele is the minor or major allele at one locus.

maize.pheno This is a `data.frame` with one column `ID` and a column `Trait` containing the measured phenotypic trait. The order of the genotypes is the same as the order of rows in `maize.geno`.

maize.ped This `data.frame` contains the pedigree information of 1301 genotypes.

maize.marker.pos This `data.frame` contains additional information for the SNP markers. The first column `pos` gives the position of the marker on the chromosome in cM. The second column `chr` specifies the chromosome (linkage group) the marker belongs to. The order of the markers is the same as the order of columns in `maize.geno`.

2 Marker data

In the first step the marker data has to be coded in a way that it could be used for the construction of genomic relationship matrices. For biallelic marker data, the minor allele should be coded as 2 and the major allele as 0. This task is done by the function `codeGeno`. If no missing values are present and all markers should be used in the following analyses, this function does the simply the recoding of the alleles. For the `maize` data, use

```
> marker <- codeGeno(maize.geno[, -1])
```

to obtain a object `marker` which contains the recoded marker data. Note that the first column is not used because it contains the `ID`. Now, the minor allele frequencies are easily obtained by dividing the column means of `marker` by 2. A histogram of the minor allele frequencies is shown in Figure 1.

In experimental data usually missing values occur in genotypic data due to different reasons (i.e. heterozygous genotypes for one locus in homozygous lines). The function `codeGeno` provides the possibility to impute missing values by chance or according to family structure by the following rules:

with population structure Suppose an observation i is missing (NA) for a marker j in population k . If marker j is fixed in population k , the imputed value will be the fixed allele. If marker j is segregating for the population k , the value is 0 with probability 0.5 and 2 with probability 0.5.

without population structure The missing values for a marker j are sampled from the allele distribution of marker j .

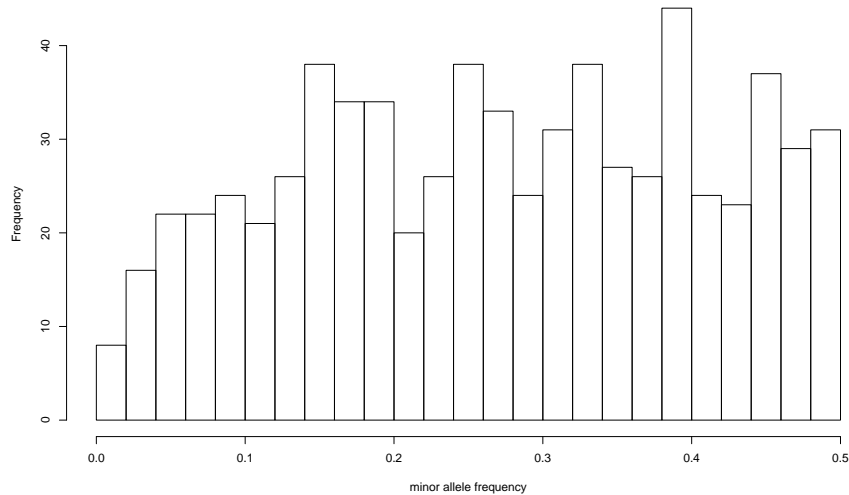


Figure 1: Histogramm of the minor allele frequency of the 696 SNP markers in **maize** data.

For illustration, 200 entries of the marker matrix are selected, the values saved and these entries are coded as NA.

```
> marker <- as.matrix(marker)
> ind1 <- sample(1:nrow(marker), 200)
> ind2 <- sample(1:ncol(marker), 200)
> original <- marker[cbind(ind1, ind2)]
> marker[cbind(ind1, ind2)] <- NA
```

The number of 1250 genotypes in the **maize** data consist of 25 *half sib* families with 50 genotypes in each family. The genotypes are ordered according to the family structure. Recoding of the marker data and imputing of the missing values is done as follows

```
> pop <- rep(1:25, each = 50)
> marker1 <- codeGeno(marker, impute = TRUE, pop)
```

```
approximative run time 6.96 seconds
```

```
...
Total number of missing values          : 200
Number of imputations by family structure : 123
Number of random imputations            : 77
Approximative fraction of correct imputations : 0.808
```

A report is printed on the screen which informs about the number of imputations performed according to family structure n_F or chance n_R . The

approximative fraction of correct imputations is $\frac{n_F+0.5n_R}{n_F+n_R}$. For the simulated data the original values are known. With the following commands the quality of the classification of the missing values is judged:

```
> imputed <- marker1[cbind(ind1, ind2)]
> (t1 <- table(original, imputed))
```

```
      imputed
original  0   2
      0 128  18
      2  24  30
```

The fraction of correct replacements is

```
> sum(diag(t1))/sum(t1)
```

```
[1] 0.79
```

In an analysis of genotypic data it is common to discard marker with a small minor allele frequency and/or many missing values. There are two additional arguments **maf** and **nmiss** for function **codeGeno**. Before recoding the data all marker with more than **nmiss**·100% missing values are discarded. After recoding the maker data only markers with a minor allele frequency > **maf** are returned by the function. By default, no markers are selected by one of both criteria.

Note that missing values in the marker data must be coded as **NA**. Instead of imputing the values **codeGeno** provides the possibility to replace the missing values by a certain value, i.e. 1 which is the expectation for the missing values. Different codings of the alleles could easily obtained with simple operations of the resulting data.

3 Pedigree

An important source of information in breeding programs is pedigree information. Especially in plant breeding pedigree is recorded over many generations. The pedigree usually consists of a list of individuals (animals or plants) of the current generation which is the suubject of analysis and their ancestors (for which usually no additional data is available). The pedigree is sorted the generation, beginning with the individuals with unknown parents. An example for an pedigree with five individuals belonging to 4 generations is given below.

Note that unknown parents are coded as "0" in **synbreed** package and generation starts with 0. In **synbreed** exists the class "pedigree", which should be used for handling pedigree information. An object of class "pedigree" consists of a **data.frame** with at least variables **ID**, **Par1**, **Par2** and

ID	Par1	Par2	gener
A	-	-	0
B	-	-	0
C	A	B	1
D	A	C	2
E	D	B	3

gener. The function `create.pedigree` creates an object of class "pedigree" for a given set of individuals and their parents. The generation can be specified by the user or otherwise is computed by the function.

Suppose we have the pedigree structure of the example. This structure is carried into `synbreed` package with the following command:

```
> id <- c("A", "B", "C", "D", "E")
> par1 <- c(0, 0, "A", "A", "D")
> par2 <- c(0, 0, "B", "C", "B")
> ped <- create.pedigree(id, par1, par2)
> ped
```

```
  ID Par1 Par2 gener
1  A    0    0     0
2  B    0    0     0
3  C    A    B     1
4  D    A    C     2
5  E    D    B     3
```

An object of class "pedigree" could be visualised with generic plotting function for S3 class "pedigree".

It is possible to simulate a pedigree structure with function `simul.pedigree`. As arguments, the number of generations to simulate and the number of individuals in each generation has to be specified. By default, random mating is assumed in each generation. As there are no further restrictions, it is possible that inbreeds could be generated when parent 1 equals parent 2. To simulate a pedigree with 6 generations and 4, 6, 7, 9, 10 and 10 individuals in each generation, use

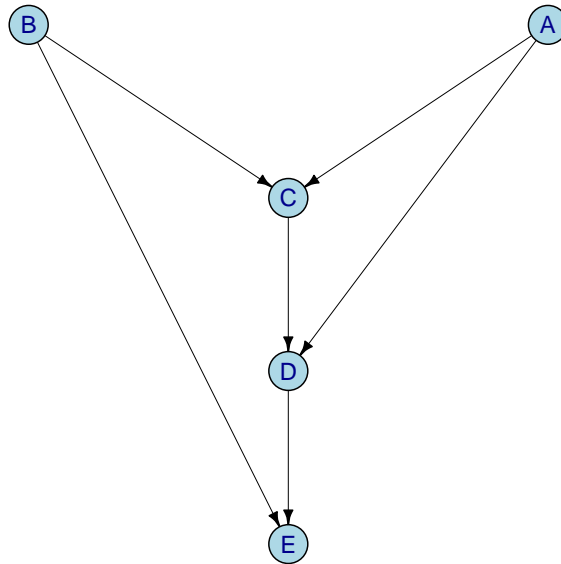
```
> set.seed(123)
> ped.simul <- simul.pedigree(gener = 6, ids = c(4, 6,
+       7, 9, 10, 10))
```

The resulting pedigree is visualized in Figure 2.

4 Relationship matrices

Pedigree information is usually used to set up relationship matrices of a number of individuals. The relationship is constructed by the *expected* fraction of alleles that are identical by descent (IBD) between relatives. Another

```
> plot(ped)
```



possibility to set up a relationship is the use of marker data to compute the genomic relationship. This gives the *observed* relationship of individuals.

4.1 Based on Pedigree

The computation of the pedigree based relationship in **synbreed** starts with the gametic relationship. A gamet is the genetic unit that an individual passes to its offspring. Thus the genetic value of an individual at one locus consists of two allelels. Suppose there is an individual C with the parents A and B. Individual C has to allelels C1 and C2. The source of allele C1 is Parent A, thus allele C1 could either equal A1 or A2. Allele C2 was inherited of parent B, thus it could be B1 or B2. To compute the gametic relationship start with an expanded table with two allelels for each individual.

This table is converted into the gametic relationship **G** matrix which is of order $2n$, if the number of individuals is n . The diagonal values of **G** are always 1. The off-diagonal values give the probabily that the two allelels are identical by descent (IBD). If the parents are unknown, it is assumed

```
> plot(ped.simul)
```

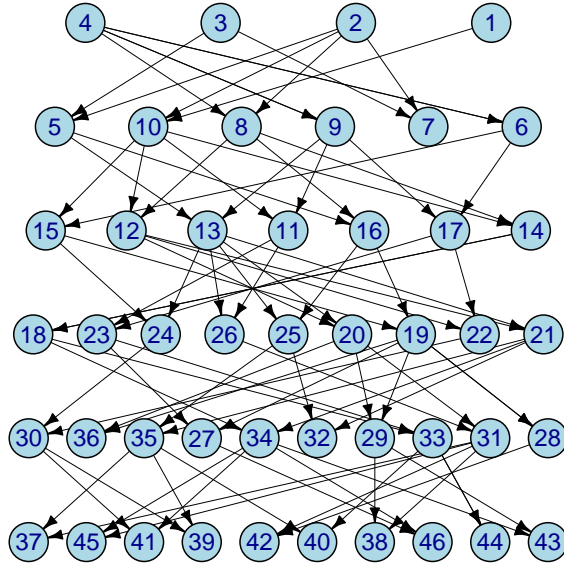


Figure 2: Simulated pedigree structure

that they are progeny of a random mating population. In this case the off-diagonals are zero. The other values of the gametic relations are filled in rowwise. The four values which describe the relationship of progeny A with parent C are computed as follows

$$(A1, C1) = 0.5 \cdot [(A1, A1) + (A1, A2)]$$

$$(A1, C2) = 0.5 \cdot [(A1, B1) + (A1, B2)]$$

$$(A2, C1) = 0.5 \cdot [(A2, A1) + (A2, A2)]$$

$$(A2, C2) = 0.5 \cdot [(A2, B1) + (A2, B2)]$$

The gametic relationship for a given pedigree is obtained as follows

```
> G <- kinship(ped, ret = "gam")
> G
```

ID	Allele	Par1	Par2
A	A1	-	-
A	A2	-	-
B	B1	-	-
B	B1	-	-
C	C1	A1	A2
C	C2	B1	B2
D	D1	A1	A2
D	D2	C1	C2
E	E1	D1	D2
E	E2	B1	B2

```

      A_1  A_2  B_1  B_2 C_1  C_2  D_1  D_2  E_1  E_2
A_1 1.000 0.000 0.000 0.000 0.5 0.00 0.500 0.250 0.375 0.000
A_2 0.000 1.000 0.000 0.000 0.5 0.00 0.500 0.250 0.375 0.000
B_1 0.000 0.000 1.000 0.000 0.0 0.50 0.000 0.250 0.125 0.500
B_2 0.000 0.000 0.000 1.000 0.0 0.50 0.000 0.250 0.125 0.500
C_1 0.500 0.500 0.000 0.000 1.0 0.00 0.500 0.500 0.500 0.000
C_2 0.000 0.000 0.500 0.500 0.0 1.00 0.000 0.500 0.250 0.500
D_1 0.500 0.500 0.000 0.000 0.5 0.00 1.000 0.250 0.625 0.000
D_2 0.250 0.250 0.250 0.250 0.5 0.50 0.250 1.000 0.625 0.250
E_1 0.375 0.375 0.125 0.125 0.5 0.25 0.625 0.625 1.000 0.125
E_2 0.000 0.000 0.500 0.500 0.0 0.50 0.000 0.250 0.125 1.000
attr(,"class")
[1] "relationshipMatrix"
```

The resulting object **G** is of class "relationshipMatrix" which is the general class for all kinds of relationship matrices (gametic relationship, additive and dominance relationship, kinship). An object of class "relationshipMatrix" is basically a symmetric matrix containing the relationship coefficient of two individuals. Note that the entry of allele 1 and allele 2 of an individual i equals his inbreeding coefficient F_i . For example, the inbreeding coefficient of individual D is

```

> G["D_1", "D_2"]

[1] 0.25
attr(,"class")
[1] "relationshipMatrix"
```

which is nonzero because individuals A and C which are the parents of D are relatives. Once the gametic relationship is computed, it could be converted in the additive numerator relationship matrix **A** or the dominance relationship matrix **D**. The additive relationship between the individuals A and B is given by

$$0.5 \cdot (G[A1, B1] + G[A1, B2] + G[A2, B1] + G[A2, B2]),$$

where $G[.,.]$ denotes the corresponding value of the gametic relationship matrix \mathbf{G} . The additive numerator relationship matrix describes the relationship between individuals and is of order n . It is typically used in the animal model

$$y_i = \mu + a_i + e_i,$$

where a_i is the additive genetic effect and $\mathbf{a} \sim N(\mathbf{0}, \mathbf{A}\sigma_{\mathbf{a}}^2)$ is assumed. The additive numerator relationship matrix for a given pedigree is obtained as follows

```
> A <- kinship(ped, ret = "add")
> A
```

```
      A      B      C      D      E
A 1.000 0.000 0.500 0.75 0.375
B 0.000 1.000 0.500 0.25 0.625
C 0.500 0.500 1.000 0.75 0.625
D 0.750 0.250 0.750 1.25 0.750
E 0.375 0.625 0.625 0.75 1.125
attr(,"class")
[1] "relationshipMatrix"
```

Note that the diagonals of \mathbf{A} are $1 + F_i$. Sometimes the kinship matrix is required which is half of the additive numerator relationship matrix. It is obtained by

```
> K <- kinship(ped, ret = "kin")
```

Additionally it is possible to derive the dominance relationship matrix \mathbf{D} out of \mathbf{G} . The dominance is needed in the non-additive animal model

$$y = \mu + a_i + a_j + d_{ij} + e_i,$$

where a_i and a_j are the additive genetic effects of alleles i and j and d_{ij} is the dominance (interaction) effect of alleles i and j with $\mathbf{d} \sim N(\mathbf{0}, \mathbf{D}\sigma_{\mathbf{d}}^2)$. The dominance relationship matrix for the example is obtained as

```
> D <- kinship(ped, ret = "dom")
> D
```

```
      A      B      C      D      E
A 1.00 0.000 0.00 0.25000 0.000000
B 0.00 1.000 0.00 0.00000 0.125000
C 0.00 0.000 1.00 0.25000 0.250000
D 0.25 0.000 0.25 1.06250 0.156250
E 0.00 0.125 0.25 0.15625 1.015625
attr(,"class")
[1] "relationshipMatrix"
```

Higher order interactions in the non-additive animal model as additive-additive, additive-dominance or dominance-dominance variance-covariance matrices can be computed as

```
> (AA <- A * A)
```

```

      A      B      C      D      E
A 1.000000 0.000000 0.250000 0.5625 0.140625
B 0.000000 1.000000 0.250000 0.0625 0.390625
C 0.250000 0.250000 1.000000 0.5625 0.390625
D 0.562500 0.062500 0.562500 1.5625 0.562500
E 0.140625 0.390625 0.390625 0.5625 1.265625
attr("class")
[1] "relationshipMatrix"
```

```
> (AD <- A * D)
```

```

      A      B      C      D      E
A 1.0000 0.000000 0.00000 0.1875000 0.0000000
B 0.0000 1.000000 0.00000 0.0000000 0.0781250
C 0.0000 0.000000 1.00000 0.1875000 0.1562500
D 0.1875 0.000000 0.18750 1.3281250 0.1171875
E 0.0000 0.078125 0.15625 0.1171875 1.1425781
attr("class")
[1] "relationshipMatrix"
```

```
> (DD <- D * D)
```

```

      A      B      C      D      E
A 1.0000 0.000000 0.0000 0.06250000 0.00000000
B 0.0000 1.000000 0.0000 0.00000000 0.01562500
C 0.0000 0.000000 1.0000 0.06250000 0.06250000
D 0.0625 0.000000 0.0625 1.12890625 0.02441406
E 0.0000 0.015625 0.0625 0.02441406 1.03149414
attr("class")
[1] "relationshipMatrix"
```

4.2 Based on marker data

The relationship matrix based on marker data or genomic relationship matrix data represents the true relationship between relatives more precise than the numerator relationship based on pedigree, as it takes into account that relationship may deviate from the expected average relationship due

to Mendelian sampling effect. Two methods for the construction of a relationship matrix based on marker data are implemented in the **synbreed** package: genomic relationship according to vanRaden (vanRaden, 2008) and according to Roger's distance.

For vanRaden, the SNP genotypes are coded as the number of copies of one of the SNP alleles, i.e., 0, 1 or 2 (any linear transformations of these values are valid too). Thus the marker data could be the result of a call of **codeGeno** when imputing for the missing values was performed or the missing values were replaced with the value 1. The genomic relationship matrix according to vanRaden is computed as

$$\frac{\mathbf{Z}\mathbf{Z}'}{2\sum_{i=1}^p p_i(1-p_i)}, \quad (1)$$

where $\mathbf{Z} = \mathbf{M} - \mathbf{P}$ and \mathbf{M} is the marker matrix and \mathbf{P} contains the allele frequencies multiplied by 2. p_i is the allele frequency of marker i . As an example we look at the marker data of 6 individuals genotyped with 8 SNP markers. Let

$$\mathbf{M} = \begin{pmatrix} 2 & 0 & 0 & 2 & 2 & 0 & 0 & 0 \\ 2 & 0 & 2 & 2 & 2 & 0 & 2 & 2 \\ 2 & 0 & 2 & 2 & 0 & 0 & 2 & 0 \\ 0 & 0 & 2 & 2 & 0 & 0 & 2 & 0 \\ 0 & 0 & 2 & 0 & 0 & 0 & 2 & 0 \\ 2 & 2 & 2 & 2 & 0 & 0 & 0 & 2 \end{pmatrix},$$

then it holds that

$$\mathbf{P} = \begin{pmatrix} 1.33 & 0.33 & 1.67 & 1.67 & 0.67 & 0 & 1.33 & 0.67 \\ 1.33 & 0.33 & 1.67 & 1.67 & 0.67 & 0 & 1.33 & 0.67 \\ 1.33 & 0.33 & 1.67 & 1.67 & 0.67 & 0 & 1.33 & 0.67 \\ 1.33 & 0.33 & 1.67 & 1.67 & 0.67 & 0 & 1.33 & 0.67 \\ 1.33 & 0.33 & 1.67 & 1.67 & 0.67 & 0 & 1.33 & 0.67 \\ 1.33 & 0.33 & 1.67 & 1.67 & 0.67 & 0 & 1.33 & 0.67 \end{pmatrix}$$

$$\mathbf{Z} = \begin{pmatrix} 0.67 & -0.33 & -1.67 & 0.33 & 1.33 & 0.00 & -1.33 & -0.67 \\ 0.67 & -0.33 & 0.33 & 0.33 & 1.33 & 0.00 & 0.67 & 1.33 \\ 0.67 & -0.33 & 0.33 & 0.33 & -0.67 & 0.00 & 0.67 & -0.67 \\ -1.33 & -0.33 & 0.33 & 0.33 & -0.67 & 0.00 & 0.67 & -0.67 \\ -1.33 & -0.33 & 0.33 & -1.67 & -0.67 & 0.00 & 0.67 & -0.67 \\ 0.67 & 1.67 & 0.33 & 0.33 & -0.67 & 0.00 & -1.33 & 1.33 \end{pmatrix}$$

and

$$\mathbf{Z}\mathbf{Z}' = \begin{pmatrix} 7.44 & 0.11 & -1.22 & -2.56 & -3.22 & -0.56 \\ 0.11 & 4.78 & -0.56 & -1.89 & -2.56 & 0.11 \\ -1.22 & -0.56 & 2.11 & 0.78 & 0.11 & -1.22 \\ -2.56 & -1.89 & 0.78 & 3.44 & 2.78 & -2.56 \\ -3.22 & -2.56 & 0.11 & 2.78 & 6.11 & -3.22 \\ -0.56 & 0.11 & -1.22 & -2.56 & -3.22 & 7.44 \end{pmatrix}$$

with the denominator $2 \sum_{i=1}^p p_i(1 - p_i) = 2.611$. To compute the genomic relationship according to vanRaden, matrix **M** is passed to the function **vanRaden**

```
> M <- matrix(data = c(2, 0, 0, 2, 2, 0, 0, 0, 2, 0, 2,
+      2, 2, 0, 2, 2, 2, 0, 2, 2, 0, 0, 2, 0, 0, 0, 2, 2,
+      0, 0, 2, 0, 0, 0, 2, 0, 0, 2, 0, 2, 2, 2, 2, 0,
+      0, 0, 2), nrow = 6, byrow = TRUE)
> vR <- vanRaden(M)
> vR
```

```
      [,1]      [,2]      [,3]      [,4]      [,5]
[1,] 2.85106383 0.04255319 -0.46808511 -0.9787234 -1.23404255
[2,] 0.04255319 1.82978723 -0.21276596 -0.7234043 -0.97872340
[3,] -0.46808511 -0.21276596 0.80851064 0.2978723 0.04255319
[4,] -0.97872340 -0.72340426 0.29787234 1.3191489 1.06382979
[5,] -1.23404255 -0.97872340 0.04255319 1.0638298 2.34042553
[6,] -0.21276596 0.04255319 -0.46808511 -0.9787234 -1.23404255
      [,6]
[1,] -0.21276596
[2,] 0.04255319
[3,] -0.46808511
[4,] -0.97872340
[5,] -1.23404255
[6,] 2.85106383
attr(,"class")
[1] "relationshipMatrix"
```

Note the object **vR** is again of class "relationshipMatrix".

Another possibility is to compute the genomic relationship matrix according to Roger's distance (Rogers, 1972). Roger's distance is computed as

$$d = \frac{1}{p} \sum_{i=1}^p \sqrt{1/2 \sum_{j=1}^{n_i} (p_{ij} - q_{ij})^2} \quad (2)$$

where p is the number of markers and n_i is the number of alleles for marker i . Let p_{ij} and q_{ij} denote the allele frequencies of allele j for marker i respectively. Note that marker data should be coded -1 and 1 for homozygous genotypes and 0 for heterozygous. If marker data is code $0/1/2$ data is transformed automatically when function **rogers** is used, which computes Roger's distance.

4.3 Doubled haploid (DH)- lines

In plant breeding doubled haploid lines are common. DH lines are fully inbred and thus have an inbreeding coefficient of 1. This has to be taken into account, when the relationship matrix in a pedigree with DH lines is computed. As an example the **maize** data is taken.

```
> data(maize)
> head(maize.ped)
```

```
   ID Par1 Par2 DH
1  1    0    0  1
2  2    0    0  1
3  3    0    0  1
4  4    0    0  1
5  5    0    0  1
6  6    0    0  1
```

First, the additive numerator relationship matrix is computed. There are 1276 DH lines and 25 non DH lines in the pedigree. For DH lines, it is necessary to set the inbreeding coefficient on 1. An argument `DH` is available for function `kinship` where for each individual in the pedigree it specifies whether this is a DH line or not. This information is available for the `maize` data. To obtain the additive numerator relationship matrix of the first 100 genotypes, use

```
> ped.maize <- create.pedigree(maize.ped$ID, maize.ped$Par1,
+   maize.ped$Par2)
> A.maize100 <- kinship(ped.maize[1:100, ], DH = maize.ped$DH[1:100],
+   ret = "add")
```

4.4 Visualisation of relationship matrices

As in most cases a relationship matrix is too big to print it on the screen. Thus there are two possibilities for visualisation of an object of class "relationshipMatrix" in `synbreed` package. A `summary` method is defined which gives the important characteristics of a relationship matrix. Use

```
> summary(A.maize100)
```

```
Dimension      : 100 x 100
Rank           : 75
Range          : 0 -- 2
# of unique values: 5
```

to get the summary for the pedigree based additive relationship matrix of the `maize` data. Another possibility is the `plot` method which could be applied to an object of class "relationshipMatrix". This gives a heatmap of the entries of the relationship matrix

Note that objects of class "relationshipMatrix" can be written to input files appropriate for `WOMBAT` (Meyer, 2006) or `ASReml` (Gilmour et al., 2000).

5 Acknowledgements

This research was funded by the German Federal Ministry of Education and Research (BMBF) within the AgroClustEr *Synbreed – Synergistic plant and animal breeding*.

References

- Gilmour, A., BR, C., SJ, W., and R, T. (2000). ASREML. program user manual. *NSW Agriculture, Orange Agricultural Institute, Forest Road, Orange, Australia*.
- Meyer, K. (2006). Wombat - a tool for mixed model analyses in quantitative genetics by reml. *J. Zhejiang Uni SCIENCE*, pages 815–821.
- Rogers, J. (1972). Measures of genetic similarity and genetic distance. In *Studies in genetics VII*, volume 7213, page 145–153. Univ. of Texas, Austin.
- vanRaden, P. (2008). Efficient methods to compute genomic predictions. *Journal of Dairy Science*, 91:4414–4423.