

# Quantitative genetics using the sommer package

*Giovanny Covarrubias-Pazaran*

2018-11-26

The sommer package was developed to provide R users a powerful and reliable multivariate mixed model solver for different genetic and non-genetic analysis in diploid and polyploid organisms. This package allows the user to estimate variance components for a mixed model with the advantage of specifying the variance-covariance structure of the random effects, specify heterogeneous variances, and obtain other parameters such as BLUPs, BLUEs, residuals, fitted values, variances for fixed and random effects, etc. The core algorithms of the package are coded in C++ using the Armadillo library to optimize dense matrix operations common in the direct-inversion algorithms.

The package is focused on problems of the type  $p > n$  related to genomic prediction (hybrid prediction & genomic selection) and GWAS analysis, although any general mixed model can be fitted as well. The package provides kernels to estimate additive (**A.mat**), dominance (**D.mat**), and epistatic (**E.mat**) relationship matrices that have been shown to increase prediction accuracy under certain scenarios or simply to estimate the variance components of such. The package provides flexibility to fit other genetic models such as full and half diallel models as well.

Vignettes aim to provide several examples in how to use the sommer package under different scenarios. We will spend the rest of the space providing examples for:

- 1) Heritability ( $h^2$ ) calculation
- 2) Specifying heterogeneous variances in mixed models
- 3) Using the pin calculator
- 4) Half and full diallel designs (using the overlay)
- 5) Genomic selection (predicting mendelian sampling)
  - GBLUP
  - rrBLUP
- 6) Single cross prediction (hybrid prediction)
- 7) Spatial modeling (using the 2-dimensional splines)
- 8) Multivariate genetic models and genetic correlations
- 9) Final remarks

## Background

The core of the package are the **mmer2** (formula-based) and **mmer** (matrix-based) functions which solve the mixed model equations. The functions are an interface to call the **NR** Direct-Inversion Newton-Raphson (Tunnicliffe 1989; Gilmour et al. 1995; Lee et al. 2016) or the **EMMA** efficient mixed model association algorithm (Kang et al. 2008). Since version 2.0 sommer can handle multivariate models. Following Maier et al. (2015), the multivariate (and by extension the univariate) mixed model implemented has the form:

$$y_1 = X_1\beta_1 + Z_1u_1 + \epsilon_1 \quad y_2 = X_2\beta_2 + Z_2u_2 + \epsilon_2 \quad \dots \quad y_i = X_i\beta_i + Z_iu_i + \epsilon_i$$

where  $y_i$  is a vector of trait phenotypes,  $\beta_i$  is a vector of fixed effects,  $u_i$  is a vector of random effects for individuals and  $\epsilon_i$  are residuals for trait 'i' ( $i = 1, \dots, t$ ). The random effects ( $u_1 \dots u_t$  and  $\epsilon_i$ ) are assumed to be normally distributed with mean zero.  $X$  and  $Z$  are incidence matrices for fixed and random effects respectively. The distribution of the multivariate response and the phenotypic variance covariance ( $V$ ) are:

$$Y = X\beta + ZU + \epsilon$$

$$Y \sim MVN(X\beta, V)$$

$$\mathbf{Y} = \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_t \end{bmatrix}$$

$$\mathbf{X} = \begin{bmatrix} X_1 & \dots & \dots \\ \vdots & \ddots & \vdots \\ \dots & \dots & X_t \end{bmatrix}$$

$$\mathbf{V} = \begin{bmatrix} Z_1 K \sigma_{g_1}^2 Z'_1 + H \sigma_{\epsilon_1}^2 & \dots & Z_1 K \sigma_{g_{1,t}} Z'_t + H \sigma_{\epsilon_{1,t}}^2 \\ \vdots & \ddots & \vdots \\ Z_1 K \sigma_{g_{1,t}} Z'_t + H \sigma_{\epsilon_{1,t}}^2 & \dots & Z_t K \sigma_{g_t}^2 Z'_t + H \sigma_{\epsilon_t}^2 \end{bmatrix}$$

where K is the relationship or covariance matrix for the kth random effect ( $u=1,\dots,k$ ), and R=I is an identity matrix for the residual term. The terms  $\sigma_{g_i}^2$  and  $\sigma_{\epsilon_i}^2$  denote the genetic (or any of the kth random terms) and residual variance of trait 'i', respectively and  $\sigma_{g_{ij}}$  and  $\sigma_{\epsilon_{ij}}$  the genetic (or any of the kth random terms) and residual covariance between traits 'i' and 'j' ( $i=1,\dots,t$ , and  $j=1,\dots,t$ ). The algorithm implemented optimizes the log likelihood:

$$\log L = 1/2 * \ln(|V|) + \ln(X'|V|X) + Y'PY$$

where  $||$  is the determinant of a matrix. And the REML estimates are updated using a Newton optimization algorithm of the form:

$$\theta^{k+1} = \theta^k + (H^k)^{-1} * \frac{dL}{d\sigma_i^2} | \theta^k$$

Where,  $\theta$  is the vector of variance components for random effects and covariance components among traits,  $H^{-1}$  is the inverse of the Hessian matrix of second derivatives for the kth cycle,  $\frac{dL}{d\sigma_i^2}$  is the vector of first derivatives of the likelihood with respect to the variance-covariance components. The Eigen decomposition of the relationship matrix proposed by Lee and Van Der Werf (2016) was included in the Newton-Raphson algorithm to improve time efficiency. Additionally, the popular pin function to estimate standard errors for linear combinations of variance components (i.e. heritabilities and genetic correlations) was added to the package as well.

The function `mmer` takes the Zs and Ks for each random effect and construct the neccesary structure inside and estimates the variance components by ML/REML using any of the 4 methods available in sommer. The `mmer2` function is enabled to work in a model-based fashion so user don't have to build the Z's and K matrices. Please refer to the canonical papers listed in the Literature section to check how the algorithms work. We have tested widely the methods to make sure they provide the same solution when the likelihood behaves well but for complex problems they might lead to slightly different answers. If you have any concern please contact me at [cova\\_ruber@live.com.mx](mailto:cova_ruber@live.com.mx).

In the following section we will go in detail over several examples on how to use mixed models in univariate and multivariate case and their use in quantitative genetics.

## 1) Marker and non-marker based heritability calculation

The heritability is one of the most popular parameters among the breeding and genetics community because of the insight that provides in the inheritance of the trait. The heritability is usually estimated as narrow sense ( $h^2$ ; only additive variance in the numerator  $\sigma_A^2$ ), and broad sense ( $H^2$ ; all genetic variance in the numerator  $\sigma_G^2$ ).

In a classical breeding experiment with no molecular markers, special designs are performed to estimate and dissect the additive ( $\sigma_A^2$ ) and non-additive (i.e. dominance  $\sigma_D^2$ ) variance along with environmental variability.

Designs such as generation analysis, North Carolina designs are used to dissect  $\sigma_A^2$  and  $\sigma_D^2$  to estimate the narrow sense heritability ( $h^2$ ). When no special design is available we can still dissect the genetic variance ( $\sigma_G^2$ ) and estimate the broad sense heritability. In this first example we will show the broad sense estimation which doesn't use covariance structures for the genotypic effect (i.e. genomic or additive relationship matrices). For big models with no covariance structures, sommer's direct inversion is a bad idea to use but we will show anyways how to do it, but keep in mind that for very sparse models we recommend using the lmer function from the lme4 package or any other package using MME-based algorithms (i.e. asreml-R).

The following dataset has 41 potato lines evaluated in 5 locations across 3 years in an RCBD design. We show how to fit the model and extract the variance components to calculate the  $h^2$ .

```
library(sommer)
data(DT_example)
head(DT)

##           Name Env Loc Year   Block Yield    Weight
## 33 Manistee(MSL292-A) CA.2013 CA.2013 CA.2013.1 4 -1.904711
## 65             C002024-9W CA.2013 CA.2013 CA.2013.1 5 -1.446958
## 66 Manistee(MSL292-A) CA.2013 CA.2013 CA.2013.2 5 -1.516271
## 67             MSL007-B CA.2011 CA.2011 CA.2011.2 5 -1.435510
## 68             MSR169-8Y CA.2013 CA.2013 CA.2013.1 5 -1.469051
## 103            AC05153-1W CA.2013 CA.2013 CA.2013.1 6 -1.307167

ans1 <- mmmer(Yield~1,
               random= ~ Name + Env + Env:Name + Env:Block,
               rcov= ~ units,
               data=DT)

## iteration  LogLik    wall   cpu(sec) restrained
## 1          -33.5019 21:28:40      0          0
## 2          -29.9296 21:28:40      0          0
## 3          -27.3285 21:28:40      0          1
## 4          -24.722   21:28:40      0          1
## 5          -24.7202 21:28:40      0          1
## 6          -24.7202 21:28:40      0          1

summary(ans1)$varcomp

##           VarComp  VarCompSE   Zratio Constraint
## Name.Yield-Yield 3.718355 1.6962316 2.1921269 Positive
## Env.Yield-Yield 12.007995 12.2729168 0.9784141 Positive
## Env:Name.Yield-Yield 5.152822 1.4926285 3.4521797 Positive
## Env:Block.Yield-Yield 0.000000 0.1156499 0.0000000 Positive
## units.Yield-Yield 4.366109 0.6572080 6.6434202 Positive

(n.env <- length(levels(DT$Env)))

## [1] 3

pin(ans1, h2 ~ V1 / (V1 + (V3/n.env) + (V5/(2*n.env)) ) )

##     Estimate       SE
## h2 0.6032719 0.1344765
```

Recently with markers becoming cheaper, thousand of markers can be run in the breeding materials. When markers are available, an special design is not neccesary to dissect the additive genetic variance. The availability of the additive, dominance and epistatic relationship matrices allow us to estimate  $\sigma_A^2$ ,  $\sigma_D^2$  and  $\sigma_I^2$ , although given that A, D and E are not orthogonal the interpretation of models that fit more than A and D become cumbersome.

Assume you have a population (even unreplicated) in the field but in addition we have genetic markers. Now we can fit the model and estimate the genomic heritability that explains a portion of the additive genetic variance (with high marker density  $\sigma_A^2 = \sigma_g^2$ )

```
data("DT_cpdata")
DT$id <- DT$id; DT$ide <- DT$id
### look at the data
A <- A.mat(GT) # additive relationship matrix
D <- D.mat(GT) # dominance relationship matrix
E <- E.mat(GT) # epistatic relationship matrix
ans.ADE <- mmmer(color~1,
                    random=~vs(id,Gu=A) + vs(idd,Gu=D),
                    rcov=~units,
                    data=DT)

## iteration    LogLik      wall    cpu(sec)   restrained
##   1       -105.511  21:28:43       1          0
##   2       -103.836  21:28:43       1          0
##   3       -103.34   21:28:43       1          0
##   4       -103.294  21:28:44       2          0
##   5       -103.293  21:28:44       2          0

(summary(ans.ADE)$varcomp)

##                               VarComp     VarCompSE   Zratio Constraint
## u:id.color-color  0.003666007 0.0012215712 3.001059 Positive
## u:idd.color-color 0.001820069 0.0007406648 2.457344 Positive
## units.color-color 0.002106117 0.0002862915 7.356547 Positive

pin(ans.ADE, h2 ~ (V1) / (V1+V3))

##      Estimate        SE
## h2 0.6351228 0.08844

pin(ans.ADE, h2 ~ (V1+V2) / (V1+V2+V3))

##      Estimate        SE
## h2 0.7225944 0.05566318
```

In the previous example we showed how to estimate the additive ( $\sigma_A^2$ ), dominance ( $\sigma_D^2$ ), and epistatic ( $\sigma_E^2$ ) variance components based on markers and estimate broad ( $H^2$ ) and narrow sense heritability ( $h^2$ ). Notice that we used the `vs()` function which indicates that the random effect inside the parenthesis (i.e. id, idd or ide) has a covariance matrix (A, D, or E), that will be specified in the Gu argument of the `vs()` function. Please DO NOT provide the inverse but the original covariance matrix.

## 2) Specifying heterogeneous variances in univariate models

Very often in multi-environment trials, the assumption that genetic variance is the same across locations may be too naive. Because of that, specifying a general genetic component and a location specific genetic variance is the way to go.

We estimate variance components for  $GCA_2$  and  $SCA$  specifying the variance structure.

```
data("DT_cornhybrids")
### fit the model
modFD <- mmmer(Yield~1,
                  random=~ vs(at(Location,c("3","4")),GCA2),
```

```

    rcov= ~ vs(ds(Location),units),
    data=DT)

## iteration LogLik     wall   cpu(sec) restrained
## 1      -200.571 21:28:45       0          0
## 2      -175.675 21:28:46       1          0
## 3      -166.325 21:28:46       1          0
## 4      -164.763 21:28:47       2          0
## 5      -164.689 21:28:47       2          0
## 6      -164.684 21:28:47       2          0
## 7      -164.684 21:28:48       3          0

summary(modFD)

## =====
## Multivariate Linear Mixed Model fit by REML
## ***** sommer 3.7 *****
## =====
##           logLik      AIC      BIC Method Converge
## Value -164.6843 331.3678 335.3592      NR      TRUE
## =====
## Variance-Covariance components:
##             VarComp VarCompSE Zratio Constraint
## 3:GCA2.Yield-Yield 62.56      53.54  1.168 Positive
## 4:GCA2.Yield-Yield 97.94      79.53  1.232 Positive
## 1:units.Yield-Yield 216.82     30.76  7.048 Positive
## 2:units.Yield-Yield 216.82     30.76  7.048 Positive
## 3:units.Yield-Yield 493.01     77.26  6.382 Positive
## 4:units.Yield-Yield 711.99    111.64  6.378 Positive
## =====
## Fixed effects:
##   Trait      Effect Estimate Std.Error t.value
## 1 Yield (Intercept) 138.1     0.9442   146.3
## =====
## Groups and observations:
##   Yield
## 3:GCA2    20
## 4:GCA2    20
## =====
## Use the '$' sign to access results and parameters

```

In the previous example we showed how the `at()` function is used in the `mmr` solver. By using the `at` function you can specify that i.e. the GCA2 has a different variance in different Locations, in this case locations 3 and 4, but also a main GCA variance. This is considered a CS + DIAG (compound symmetry + diagonal) model.

In addition, other functions can be added on top to fit models with covariance structures, i.e. the `Gu` argument from the `vs()` function to indicate a covariance matrix (A, pedigree or genomic relationship matrix)

```

data("DT_cornhybrids")
GT[1:4,1:4]

##          A258        A634        A641        A680
## A258  2.23285528 -0.3504778 -0.04756856 -0.32239362
## A634 -0.35047780  1.4529169  0.45203869 -0.02293680
## A641 -0.04756856  0.4520387  1.96940221 -0.09896791
## A680 -0.32239362 -0.0229368 -0.09896791  1.65221984

```

```

#### fit the model
modFD <- mmmer(Yield~1,
                  random=~ vs(at(Location,c("3","4")),GCA2,Gu=GT),
                  rcov= ~ vs(ds(Location),units),
                  data=DT)

## iteration    LogLik      wall    cpu(sec)   restrained
##   1     -208.145  21:28:49       1           0
##   2     -181.594  21:28:49       1           0
##   3     -169.58   21:28:50       2           0
##   4     -165.782  21:28:50       2           0
##   5     -165.279  21:28:51       3           0
##   6     -165.233  21:28:51       3           0
##   7     -165.229  21:28:52       4           0
##   8     -165.229  21:28:52       4           0

summary(modFD)

## =====
## Multivariate Linear Mixed Model fit by REML
## **** sommer 3.7 ****
## =====
##      logLik      AIC      BIC Method Converge
## Value -165.2289 332.4572 336.4487      NR      TRUE
## =====
## Variance-Covariance components:
##          VarComp VarCompSE Zratio Constraint
## 3:GCA2.Yield-Yield  26.65    26.17 1.0183 Positive
## 4:GCA2.Yield-Yield  37.56    37.84 0.9926 Positive
## 1:units.Yield-Yield 216.77   30.75 7.0490 Positive
## 2:units.Yield-Yield 216.77   30.75 7.0490 Positive
## 3:units.Yield-Yield 503.61   77.87 6.4676 Positive
## 4:units.Yield-Yield 738.80  114.15 6.4723 Positive
## =====
## Fixed effects:
##   Trait      Effect Estimate Std.Error t.value
## 1 Yield (Intercept) 138.1    0.9147    151
## =====
## Groups and observations:
##   Yield
## 3:GCA2    20
## 4:GCA2    20
## =====
## Use the '$' sign to access results and parameters

```

### 3) Using the pin calculator

Sometimes the user needs to calculate ratios or functions of specific variance-covariance components and obtain the standard error for such parameters. Examples of these are the genetic correlations, heritabilities, etc. Using the CPdata we will show how to estimate the heritability and the standard error using the pin function that uses the delta method to come up with these parameters. This can be extended for any linear combination of the variance components.

```

data("DT_cpdata")
### look at the data
A <- A.mat(GT) # additive relationship matrix
ans <- mmer(color~1,
             random=~vs(id,Gu=A),
             rcov=~units,
             data=DT)

## iteration    LogLik      wall    cpu(sec)   restrained
##    1     -110.774  21:28:53       0          0
##    2     -110.751  21:28:53       0          0
##    3     -110.742  21:28:54       1          0
##    4     -110.741  21:28:54       1          0
##    5     -110.741  21:28:54       1          0

(summary(ans.ADE)$varcomp

##                               VarComp   VarCompSE   Zratio Constraint
## u:id.color-color  0.003666007 0.0012215712 3.001059 Positive
## u:idd.color-color 0.001820069 0.0007406648 2.457344 Positive
## units.color-color 0.002106117 0.0002862915 7.356547 Positive

pin(ans, h2 ~ (V1) / (V1+V2))

##      Estimate        SE
## h2 0.6512863 0.06109601

```

The same can be used for multivariate models. Please check the documentation of the `pin` function to see more examples.

#### 4) Half and full diallel designs (use of the overlay)

When breeders are looking for the best single cross combinations, diallel designs have been by far the most used design in crops like maize. There are 4 types of diallel designs depending if reciprocate and self cross (omission of parents) are performed (full diallel with parents  $n^2$ ; full diallel without parents  $n(n-1)$ ; half diallel with parents  $1/2 * n(n+1)$ ; half diallel without parents  $1/2 * n(n-1)$ ). In this example we will show a full diallel design (reciprocate crosses are performed) and half diallel designs (only one of the directions is performed).

In the first data set we show a full diallel among 40 lines from 2 heterotic groups, 20 in each. Therefore 400 possible hybrids are possible. We have phenotypic data for 100 of them across 4 locations. We use the data available to fit a model of the form:

$$y = X\beta + Zu_1 + Zu_2 + Zu_S + \epsilon$$

We estimate variance components for  $GCA_1$ ,  $GCA_2$  and  $SCA$  and use them to estimate heritability. Additionally BLUPs for GCA and SCA effects can be used to predict crosses.

```

data("DT_cornhybrids")

modFD <- mmer(Yield~Location,
                random=~GCA1+GCA2+SCA,
                rcov=~units,
                data=DT)

## iteration    LogLik      wall    cpu(sec)   restrained
##    1     -162.189  21:28:55       0          0

```

```

##   2    -149.491  21:28:55      0      0
##   3    -138.221  21:28:56      1      1
##   4    -132.793  21:28:56      1      1
##   5    -132.628  21:28:57      2      1
##   6    -132.597  21:28:57      2      1
##   7    -132.59    21:28:57      2      1
##   8    -132.589  21:28:58      3      1
##   9    -132.589  21:28:58      3      1

(suma <- summary(modFD)$varcomp)

##                               VarComp VarCompSE     Zratio Constraint
## GCA1.Yield-Yield    0.000000  16.50320  0.0000000  Positive
## GCA2.Yield-Yield    7.416668  18.94490  0.3914864  Positive
## SCA.Yield-Yield   187.556634  41.59316  4.5093139  Positive
## units.Yield-Yield 221.142456  18.14715 12.1860689  Positive

Vgca <- sum(suma[1:2,1])
Vsca <- suma[3,1]
Ve <- suma[4,1]
Va = 4*Vgca
Vd = 4*Vsca
Vg <- Va + Vd
(H2 <- Vg / (Vg + (Ve)) )

## [1] 0.7790863

(h2 <- Va / (Vg + (Ve)) )

## [1] 0.02963598

```

Don't worry too much about the small h2 value, the data was simulated to be mainly dominance variance, therefore the Va was simulated extremely small leading to such value of narrow sense h2.

In this second data set we show a small half diallel with 7 parents crossed in one direction.  $n(n-1)/2$  crosses are possible  $7(6)/2 = 21$  unique crosses. Parents appear as males or females indistinctly. Each with two replications in a CRD. For a half diallel design a single GCA variance component for both males and females can be estimated and an SCA as well ( $\sigma_G^2 CA$  and  $\sigma_S^2 CA$  respectively), and BLUPs for GCA and SCA of the parents can be extracted. We would show first how to use it with the `mmer2` function using the `overlay()` function and later we will show how to do it creating customized matrices using the `overlay` and `model.matrix` functions for the GCA and SCA matrices respectively. The specific model here is:

$$y = X\beta + Zu_g + Zu_s + \epsilon$$

```

data("DT_halfdiallel")
head(DT)

##   rep geno male female      sugar
## 1   1    12     1      2 13.950509
## 2   2    12     1      2  9.756918
## 3   1    13     1      3 13.906355
## 4   2    13     1      3  9.119455
## 5   1    14     1      4  5.174483
## 6   2    14     1      4  8.452221

DT$femalef <- as.factor(DT$female)
DT$malef <- as.factor(DT$male)
DT$genof <- as.factor(DT$geno)
#### model using overlay

```

```

modh <- mmmer(sugar~1,
               random=~vs(overlay(femalef,malef))
               + genof,
               data=DT)

## iteration    LogLik      wall    cpu(sec)   restrained
##    1     -7.04379  21:28:58       0          0
##    2     -6.09505  21:28:58       0          0
##    3     -5.71831  21:28:58       0          0
##    4     -5.67487  21:28:58       0          0
##    5     -5.67441  21:28:58       0          0

summary(modh)$varcomp

##                               VarComp VarCompSE   Zratio Constraint
## u:femalef.sugar-sugar 5.508557  3.578396 1.539393 Positive
## genof.sugar-sugar     1.816367  1.364196 1.331456 Positive
## units.sugar-sugar    3.117182  0.961511 3.241962 Positive

```

Notice how the `overlay()` argument makes the overlap of incidence matrices possible making sure that male and female are joint into a single random effect.

## 5) Genomic selection

In this section we will use wheat data from CIMMYT to show how is genomic selection performed. This is the case of prediction of specific individuals within a population. It basically uses a similar model of the form:

$$y = X\beta + Zu + \epsilon$$

and takes advantage of the variance covariance matrix for the genotype effect known as the additive relationship matrix (A) and calculated using the `A.mat` function to establish connections among all individuals and predict the BLUPs for individuals that were not measured. The prediction accuracy depends on several factors such as the heritability ( $h^2$ ), training population used (TP), size of TP, etc.

```

data("DT_wheat");
colnames(DT) <- paste0("X",1:ncol(DT))
DT <- as.data.frame(DT);DT$id <- as.factor(rownames(DT))
# select environment 1
rownames(GT) <- rownames(DT)
K <- A.mat(GT) # additive relationship matrix
colnames(K) <- rownames(K) <- rownames(DT)
# GBLUP pedigree-based approach
set.seed(12345)
y.trn <- DT
vv <- sample(rownames(DT),round(nrow(DT)/5))
y.trn[vv,"X1"] <- NA
head(y.trn)

```

	X1	X2	X3	X4	id
## 775	NA	-1.72746986	-1.89028479	0.0509159	775
## 2166	-0.2527028	0.40952243	0.30938553	-1.7387588	2166
## 2167	0.3418151	-0.64862633	-0.79955921	-1.0535691	2167
## 2465	NA	0.09394919	0.57046773	0.5517574	2465
## 3881	NA	-0.28248062	1.61868192	-0.1142848	3881
## 3889	2.3360969	0.62647587	0.07353311	0.7195856	3889

```

## GBLUP
ans <- mmer(X1~1,
              random=~vs(id,Gu=K),
              rcov=~units,
              data=y.trn) # kinship based

## iteration    LogLik      wall    cpu(sec)  restrained
##   1       -211.984  21:29:0      1          0
##   2       -202.68   21:29:0      1          0
##   3       -198.262  21:29:1      2          0
##   4       -197.526  21:29:2      3          0
##   5       -197.508  21:29:2      3          0
##   6       -197.508  21:29:2      3          0

ans$U$`u:id`$X1 <- as.data.frame(ans$U$`u:id`$X1)
rownames(ans$U$`u:id`$X1) <- gsub("id","",rownames(ans$U$`u:id`$X1))
cor(ans$U$`u:id`$X1[vv,],DT[vv,"X1"], use="complete")

## [1] 0.4885724

## rrBLUP
ans2 <- mmer(X1~1,
               random=~vs(list(GT)),
               rcov=~units,
               data=y.trn) # kinship based

## iteration    LogLik      wall    cpu(sec)  restrained
##   1       -391.485  21:29:5      2          0
##   2       -259.534  21:29:5      2          0
##   3       -212.267  21:29:6      3          0
##   4       -198.261  21:29:6      3          0
##   5       -197.526  21:29:7      4          0
##   6       -197.508  21:29:7      4          0
##   7       -197.508  21:29:8      5          0

u <- GT %*% as.matrix(ans2$U$`u:GT`$X1) # BLUPs for individuals
rownames(u) <- rownames(GT)
cor(u[vv,],DT[vv,"X1"]) # same correlation

## [1] 0.4885724
# the same can be applied in multi-response models in GBLUP or rrBLUP

```

## 6) Single cross prediction

When doing prediction of single cross performance the phenotype can be dissected in three main components, the general combining abilities (GCA) and specific combining abilities (SCA). This can be expressed with the same model analyzed in the diallel experiment mentioned before:

$$y = X\beta + Zu_1 + Zu_2 + Zu_S + \epsilon$$

with:

$$u_1 \sim N(0, K_1 \sigma_u^2 1)$$

$$u_2 \sim N(0, K_2 \sigma_u^2 2)$$

$$u_s \sim N(0, K_3 \sigma_u^2 s)$$

And we can specify the K matrices. The main difference between this model and the full and half diallel designs is the fact that this model will include variance covariance structures in each of the three random effects (GCA1, GCA2 and SCA) to be able to predict the crosses that have not occurred yet. We will use the data published by Technow et al. (2015) to show how to do prediction of single crosses.

```

data("DT_technow")
# RUN THE PREDICTION MODEL
y.trn <- DT
vv1 <- which(!is.na(DT$GY))
vv2 <- sample(vv1, 100)
y.trn[vv2, "GY"] <- NA
anss2 <- mmer(GY~1,
               random=~vs(dent, Gu=Ad) + vs(flint, Gu=Af),
               rcov=~units,
               data=y.trn)

## iteration    LogLik      wall     cpu(sec)   restrained
##    1       131.469  21:29:18        8          0
##    2       140.706  21:29:25       15          0
##    3       145.859  21:29:33       23          0
##    4       147.09   21:29:40       30          0
##    5       147.178  21:29:47       37          0
##    6       147.184  21:29:55       45          0
##    7       147.184  21:30:2        52          0

summary(anss2)$varcomp

##           VarComp VarCompSE   Zratio Constraint
## u:dent.GY-GY 16.94000 2.6936588 6.288845 Positive
## u:flint.GY-GY 12.48251 2.3310159 5.354967 Positive
## units.GY-GY   16.74885 0.7660229 21.864680 Positive

zu1 <- model.matrix(~dent-1,y.trn) %*% anss2$U$`u:dent`$GY
zu2 <- model.matrix(~flint-1,y.trn) %*% anss2$U$`u:flint`$GY
u <- zu1+zu2+anss2$Beta[1,"Estimate"]
cor(u[vv2,], DT$GY[vv2])

## [1] 0.8234455

```

In the previous model we only used the GCA effects (GCA1 and GCA2) for practicality, although it's been shown that the SCA effect doesn't actually help that much in increasing prediction accuracy and increase a lot the computation intensity required since the variance covariance matrix for SCA is the kronecker product of the variance covariance matrices for the GCA effects, resulting in a 10578x10578 matrix that increases in a very intensive manner the computation required.

A model without covariance structures would show that the SCA variance component is insignificant compared to the GCA effects. This is why including the third random effect doesn't increase the prediction accuracy.

## 7) Spatial modeling (using the 2-dimensional spline)

We will use the CPdata to show the use of 2-dimensional splines for accomodating spatial effects in field experiments. In early generation variety trials the availability of seed is low, which makes the use of unreplicated design a necessity more than anything else. Experimental designs such as augmented designs and partially-replicated (p-rep) designs become every day more common these days.

In order to do a good job modeling the spatial trends happening in the field special covariance structures have been proposed to accomodate such spatial trends (i.e. autoregressive residuals; ar1). Unfortunately,

some of these covariance structures make the modeling rather unstable. More recently other research groups have proposed the use of 2-dimensional splines to overcome such issues and have a more robust modeling of the spatial terms (Lee et al. 2013; Rodríguez-Álvarez et al. 2018).

In this example we assume an unreplicated population where row and range information is available which allows us to fit a 2 dimensional spline model.

```

data("DT_cpdata")
### mimic two fields
A <- A.mat(GT)
mix <- mmmer(Yield~1,
  random=~vs(id, Gu=A) +
    vs(Rowf) +
    vs(Colf) +
    vs(sp12D(Row,Col)),
  rcov=~vs(units),
  data=DT)

## iteration      LogLik      wall      cpu(sec)      restrained
##   1      -189.212  21:30:5          1              0
##   2      -168.339  21:30:5          1              0
##   3      -154.84   21:30:6          2              0
##   4      -151.445  21:30:6          2              0
##   5      -151.225  21:30:6          2              0
##   6      -151.203  21:30:7          3              0
##   7      -151.201  21:30:7          3              0
##   8      -151.201  21:30:7          3              0

summary(mix)

## =====
##           Multivariate Linear Mixed Model fit by REML
## **** sommer 3.7 ****
## =====
##      logLik      AIC      BIC Method Converge
## Value -151.2012 304.4021 308.2937      NR      TRUE
## =====
## Variance-Covariance components:
##                               VarComp VarCompSE Zratio Constraint
## u:id.Yield-Yield      783.3     319.2 2.4540  Positive
## u:Rowf.Yield-Yield    814.9     391.0 2.0840  Positive
## u:Colf.Yield-Yield    182.2     129.6 1.4056  Positive
## u:Row.Yield-Yield     513.4     694.4 0.7393  Positive
## u:units.Yield-Yield   2922.7    294.1 9.9365  Positive
## =====
## Fixed effects:
##   Trait      Effect Estimate Std.Error t.value
## 1 Yield (Intercept) 132.1     8.792   15.03
## =====
## Groups and observations:
##      Yield
## u:id      363
## u:Rowf     13
## u:Colf     36
## u:Row     168
## =====

```

```
## Use the '$' sign to access results and parameters
```

Notice that the job is done by the `spl2D()` function that takes the Row and Col information to fit a spatial kernel.

## 8) Multivariate genetic models and genetic correlations

Sometimes is important to estimate genetic variance-covariance among traits, multi-reponse models are very useful for such task. Let see an example with 3 traits (color, Yield, and Firmness) and a single random effect (genotype; id) although multiple effects can be modeled as well. We need to use a variance covariance structure for the random effect to be able to obtain the genetic covariance among traits.

```
data("DT_cpdata")
A <- A.mat(GT)
ans.m <- mmmer(cbind(Yield,color)^~1,
  random=~ vs(id, Gw=A)
  + vs(Rowf,Gtc=diag(2))
  + vs(Colf,Gtc=diag(2)),
  rcov=~ vs(units),
  data=DT)
```

## iteration	LogLik	wall	cpu(sec)	restrained
## 1	-407.537	21:30:13	5	0
## 2	-347.164	21:30:17	9	0
## 3	-283.864	21:30:22	14	0
## 4	-255.884	21:30:26	18	0
## 5	-253.53	21:30:31	23	0
## 6	-253.305	21:30:35	27	0
## 7	-253.28	21:30:39	31	0
## 8	-253.277	21:30:44	36	0
## 9	-253.277	21:30:48	40	0

Now you can extract the BLUPs using the ‘randef’ function or simple accesing with the ‘\$’ sign and pick ‘u.hat’. Also, genetic correlations and heritabilities can be calculated easily.

```
cov2cor(ans.m$sigma$`u:id`)
```

```
##          Yield      color
## Yield  1.0000000  0.1231828
## color   0.1231828  1.0000000
```

## 9) Final remarks

Keep in mind that sommer uses direct inversion (DI) algorithm which can be very slow for large datasets. The package is focused in problems of the type  $p > n$  (more random effect levels than observations) and models with dense covariance structures. For example, for experiment with dense covariance structures with low-replication (i.e. 2000 records from 1000 individuals replicated twice with a covariance structure of 1000x1000) sommer will be faster than MME-based software. Also for genomic problems with large number of random effect levels, i.e. 300 individuals ( $n$ ) with 100,000 genetic markers ( $p$ ). For highly replicated trials with small covariance structures or  $n > p$  (i.e. 2000 records from 200 individuals replicated 10 times with covariance structure of 200x200) asreml or other MME-based algorithms will be much faster and we recommend you to opt for those software.

## Literature

- Covarrubias-Pazaran G. 2016. Genome assisted prediction of quantitative traits using the R package sommer. PLoS ONE 11(6):1-15.
- Bernardo Rex. 2010. Breeding for quantitative traits in plants. Second edition. Stemma Press. 390 pp.
- Gilmour et al. 1995. Average Information REML: An efficient algorithm for variance parameter estimation in linear mixed models. Biometrics 51(4):1440-1450.
- Henderson C.R. 1975. Best Linear Unbiased Estimation and Prediction under a Selection Model. Biometrics vol. 31(2):423-447.
- Kang et al. 2008. Efficient control of population structure in model organism association mapping. Genetics 178:1709-1723.
- Lee, D.-J., Durban, M., and Eilers, P.H.C. (2013). Efficient two-dimensional smoothing with P-spline ANOVA mixed models and nested bases. Computational Statistics and Data Analysis, 61, 22 - 37.
- Lee et al. 2015. MTG2: An efficient algorithm for multivariate linear mixed model analysis based on genomic information. Cold Spring Harbor. doi: <http://dx.doi.org/10.1101/027201>.
- Maier et al. 2015. Joint analysis of psychiatric disorders increases accuracy of risk prediction for schizophrenia, bipolar disorder, and major depressive disorder. Am J Hum Genet; 96(2):283-294.
- Rodriguez-Alvarez, Maria Xose, et al. Correcting for spatial heterogeneity in plant breeding experiments with P-splines. Spatial Statistics 23 (2018): 52-71.
- Searle. 1993. Applying the EM algorithm to calculating ML and REML estimates of variance components. Paper invited for the 1993 American Statistical Association Meeting, San Francisco.
- Yu et al. 2006. A unified mixed-model method for association mapping that accounts for multiple levels of relatedness. Genetics 38:203-208.
- Abdollahi Arpanahi R, Morota G, Valente BD, Kranis A, Rosa GJM, Gianola D. 2015. Assessment of bagging GBLUP for whole genome prediction of broiler chicken traits. Journal of Animal Breeding and Genetics 132:218-228.
- Tunnicliffe W. 1989. On the use of marginal likelihood in time series model estimation. JRSS 51(1):15-27.