

Quick start for the sommer package

Giovanny Covarrubias-Pazaran

2018-03-03

The sommer package was developed to provide R users a powerful and reliable multivariate mixed model solver. The package is focused in problems of the type $p > n$ (more random effect levels than observations). This package allows the user to fit mixed models with the advantage of specifying the variance-covariance structure for the random effects, and specify heterogeneous variances, and obtain other parameters such as BLUPs, BLUEs, residuals, fitted values, variances for fixed and random effects, etc.

The purpose of this quick start guide is to show the flexibility of the package under certain common scenarios:

B1) Background on mixed models B2) Background on covariance structures

- 1) Univariate homogeneous variance models
- 2) Univariate heterogeneous variance models
- 3) Univariate unstructured variance models
- 4) Multivariate homogeneous variance models
- 5) Multivariate heterogeneous variance models
- 6) Including special functions
- 7) Spatial modeling

B1) Background on mixed models

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). From 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:

$$y_1 = X_1\beta_1 + Z_1u_1 + \epsilon_1$$

$$y_2 = X_2\beta_2 + Z_2u_2 + \epsilon_2$$

...

$$y_i = X_i\beta_i + Z_iu_i + \epsilon_i$$

where y_i is a vector of trait phenotypes, β_i is a vector of fixed effects, u_i is a vector of random effects for individuals and e_i are residuals for trait 'i' ($i = 1, \dots, t$). The random effects ($u_1 \dots u_i$ and e_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:

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

$$Y \sim \text{MVN}(X\beta, V)$$

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

$$\mathbf{X} = \begin{bmatrix} X_1 & 0 & 0 \\ \vdots & \ddots & \vdots \\ 0 & 0 & X_t \end{bmatrix}$$

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

where \mathbf{K} is the relationship or covariance matrix for the k th random effect ($u=1, \dots, k$), and $\mathbf{R}=\mathbf{I}$ 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 k th random terms) and residual variance of trait 'i', respectively and $\sigma_{g_{ij}}$ and $\sigma_{\epsilon_{ij}}$ the genetic (or any of the k th 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, θ 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 k th 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 Z s and K s for each random effect and construct the necessary 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`.

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.

B2) Background on covariance structures

One of the major strenghts of linear mixed models is the flexibility to specify variance-covariance structures at all levels. In general, mixed models can be seen as kronecker products of multiple variance-covariance stuctures. For example, a multivariate model (i.e. 2 traits) where "g" genotypes (i.e. 100 genotypes) are tested in "e" environments (i.e. 3 environments), the genotype variance-covariance can be seen as the following multiplicative model:

$$\mathbf{T} \otimes \mathbf{G} \otimes \mathbf{A}$$

where:

$$\mathbf{T} = \begin{bmatrix} \sigma_{g_{t1,t1}}^2 & \sigma_{g_{t1,t2}} \\ \sigma_{g_{t2,t1}} & \sigma_{g_{t2,t2}}^2 \end{bmatrix}$$

is the covariance structure for genotypes among traits.

$$\mathbf{G} = \begin{bmatrix} \sigma_{g_{e1,e1}}^2 & \sigma_{g_{e1,e2}} & \sigma_{g_{e1,e3}} \\ \sigma_{g_{e2,e1}} & \sigma_{g_{e2,e2}}^2 & \sigma_{g_{e2,e3}} \\ \sigma_{g_{e3,e1}} & \sigma_{g_{e3,e2}} & \sigma_{g_{e3,e3}}^2 \end{bmatrix}$$

is the genotype covariance structure among environments.

and A is the genomic, additive or any other relationship matrix.

The T and G covariance structures shown above are known as unstructured (US) covariance matrices, although this is just one example from several covariance structures that the linear mixed models enable. For example, other popular covariance structures are:

Diagonal (DIAG)

$$\mathbf{\Sigma} = \begin{bmatrix} \sigma_{g_{e1,e1}}^2 & 0 & 0 \\ \vdots & \ddots & \vdots \\ 0 & 0 & \sigma_{g_{ei,ei}}^2 \end{bmatrix}$$

Compound symmetry (CS)

$$\mathbf{\Sigma} = \begin{bmatrix} \sigma_g^2 + \sigma_{ge}^2 & \sigma_g^2 & \sigma_g^2 & \sigma_g^2 \\ \sigma_g^2 & \sigma_g^2 + \sigma_{ge}^2 & \sigma_g^2 & \sigma_g^2 \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_g^2 & \sigma_g^2 & \sigma_g^2 & \sigma_g^2 + \sigma_{ge}^2 \end{bmatrix}$$

First order autoregressive (AR1)

$$\mathbf{\Sigma} = \sigma^2 \begin{bmatrix} 1 & \rho & \rho^2 & \rho^3 \\ \rho & 1 & \rho & \rho^2 \\ \rho^2 & \rho & 1 & \rho \\ \rho^3 & \rho^2 & \rho & 1 \end{bmatrix}$$

and the already mentioned Unstructured (US)

$$\mathbf{\Sigma} = \begin{bmatrix} \sigma_{g_{e1,e1}}^2 & \sigma_{g_{e1,e2}} & \sigma_{g_{e1,e3}} \\ \vdots & \ddots & \vdots \\ \sigma_{g_{e3,e1}} & \sigma_{g_{e3,e2}} & \sigma_{g_{e3,e3}}^2 \end{bmatrix}$$

among others. `Sommer` has the capabilities to fit some of these covariance structures in the mixed model machinery. In the following section we will go over some examples on how to accommodate some structures.

1) Univariate homogeneous variance models

This type of models refer to single response models where a variable of interest (i.e. genotypes) needs to be analyzed as interacting with a 2nd random effect (i.e. environments), but you assume that across environments the genotypes have the same variance component. This is the so-called compound symmetry (CS) model.

```
library(sommer)
data(example)
head(example)
```

```
##           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 <- mmer2(Yield-Env,
              random= ~ Name + Env:Name,
              rcov= ~ units,
              data=example, silent = TRUE)
summary(ans1)
```

```
## =====
##      Multivariate Linear Mixed Model fit by REML
## ***** sommer 3.3 *****
## =====
##           logLik      AIC      BIC Method Converge
## Value -20.14537 46.29075 55.95182      MNR      TRUE
## =====
## Variance-Covariance components:
##           VarComp VarCompSE Zratio
## Name.Yield-Yield      3.682      1.6912  2.177
## Env:Name.Yield-Yield  5.173      1.4955  3.459
## units.Yield-Yield      4.366      0.6469  6.749
## =====
## Fixed effects:
##
## $Yield
##           Estimate Std. Error  t value
## (Intercept) 16.496351  0.6855051 24.064519
## EnvCA.2012  -5.776758  0.7558233 -7.643000
## EnvCA.2013  -6.380479  0.7960572 -8.015101
##
## =====
## Groups and observations:
##           Observ Groups
## Name      185      41
## Env:Name  185     123
## =====
## Use the '$' sign to access results and parameters
```

2) Univariate heterogeneous variance models

Very often in multi-environment trials, the assumption that the genetic variance or the residual 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. This requires a CS+DIAG model (also called heterogeneous CS model).

```
data(example)
head(example)
```

```
##           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 <- mmer2(Yield-Env,
              random= ~Name + at(Env):Name,
              rcov= ~ at(Env):units,
              data=example, silent = TRUE)
summary(ans1)
```

```
## =====
##      Multivariate Linear Mixed Model fit by REML
## ***** sommer 3.3 *****
## =====
##      logLik      AIC      BIC Method Converge
## Value -15.42982 36.85964 46.52071   MNR      TRUE
## =====
## Variance-Covariance components:
##      VarComp VarCompSE Zratio
## Name.Yield-Yield      2.962  1.4964  1.980
## CA.2011:Name.Yield-Yield 10.148  4.5107  2.250
## CA.2012:Name.Yield-Yield  1.879  1.8700  1.005
## CA.2013:Name.Yield-Yield  6.629  2.5027  2.649
## CA.2013:units.Yield-Yield 2.560  0.6398  4.001
## CA.2011:units.Yield-Yield 4.942  1.5246  3.242
## CA.2012:units.Yield-Yield 5.725  1.3119  4.364
## =====
## Fixed effects:
##
## $Yield
##      Estimate Std. Error  t value
## (Intercept) 16.507676  0.8268629 19.964224
## EnvCA.2012  -5.816887  0.8575779 -6.782926
## EnvCA.2013  -6.412430  0.9356441 -6.853493
##
## =====
## Groups and observations:
##      Observ Groups
## Name      185    41
## CA.2011:Name 185    41
## CA.2012:Name 185    41
## CA.2013:Name 185    41
## =====
## Use the '$' sign to access results and parameters
```

As you can see the special function `at` or `diag` can be used to indicate that there's a different variance for the genotypes in each environment. Same was done for the residual. The difference between `at` and `diag` is that the `at` function can be used to specify the levels or specific environments where the variance is different.

3) Unstructured variance models

A more relaxed assumption than the CS+DIAG model is the unstructured model (US) which assumes that among the levels of certain factor (i.e. Environments) there's a covariance structure of a second random effect (i.e. Genotypes). This can be done in sommer using the `us(.)` function:

```
data(example)
head(example)
```

```
##           Name      Env Loc Year      Block Yield      Weight
## 33  Manistee(MSL292-A) CA.2013  CA 2013  CA.2013.1      4 -1.904711
## 65           CO02024-9W 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
```

```
ans3 <- mmer2(Yield~Env,
              random=~ us(Env):Name,
              rcov=~at(Env):units,
              data=example, silent=TRUE)
summary(ans3)
```

```
## =====
##           Multivariate Linear Mixed Model fit by REML
## ***** sommer 3.3 *****
## =====
##           logLik      AIC      BIC Method Converge
## Value -11.4997 28.99939 38.66046 MNR TRUE
## =====
## Variance-Covariance components:
##           VarComp VarCompSE Zratio
## Env:Name!Env.CA.2011:CA.2011.Yield-Yield 15.6650 5.4206 2.8899
## Env:Name!Env.CA.2012:CA.2011.Yield-Yield 6.1109 2.4858 2.4583
## Env:Name!Env.CA.2013:CA.2011.Yield-Yield 6.3841 3.0659 2.0823
## Env:Name!Env.CA.2012:CA.2012.Yield-Yield 4.5309 1.8217 2.4872
## Env:Name!Env.CA.2013:CA.2012.Yield-Yield 0.3916 1.5244 0.2569
## Env:Name!Env.CA.2013:CA.2013.Yield-Yield 8.5978 2.4844 3.4607
## CA.2013:units.Yield-Yield 2.5570 0.6391 4.0008
## CA.2011:units.Yield-Yield 4.9699 1.5323 3.2434
## CA.2012:units.Yield-Yield 5.6723 1.3001 4.3631
## =====
## Fixed effects:
##
## $Yield
##           Estimate Std. Error t value
## (Intercept) 16.331260 0.8137093 20.070141
## EnvCA.2012 -5.695867 0.7403739 -7.693229
## EnvCA.2013 -6.271133 0.8191001 -7.656125
##
## =====
## Groups and observations:
##           Observ Groups
## Env:Name!Env.CA.2011:CA.2011 185 41
## Env:Name!Env.CA.2012:CA.2011 185 82
```

```
## Env:Name!Env.CA.2013:CA.2011      185      82
## Env:Name!Env.CA.2012:CA.2012      185      41
## Env:Name!Env.CA.2013:CA.2012      185      82
## Env:Name!Env.CA.2013:CA.2013      185      41
## =====
## Use the '$' sign to access results and parameters
```

As can be seen the `us(Env)` indicates that the genotypes (Name) can have a covariance structure among environments (Env).

4) Multivariate homogeneous variance models

Currently there's a great push for multi-response models. This is motivated by the correlation that certain variables hide and that could benefit in the prediction perspective. In sommer to specify multivariate models the response requires the use of the `cbind()` function in the response, and the `us(trait)`, `diag(trait)`, or `at(trait)` functions in the random part of the model.

```
data(example)
head(example)
```

```
##           Name      Env Loc Year      Block Yield      Weight
## 33  Manistee(MSL292-A) CA.2013  CA 2013  CA.2013.1      4 -1.904711
## 65           CO02024-9W 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 <- mmer2(cbind(Yield, Weight) ~ Env,
              random= ~ us(trait):Name + us(trait):Env:Name,
              rcov= ~ us(trait):units,
              data=example, silent = TRUE)
summary(ans1)
```

```
## =====
##      Multivariate Linear Mixed Model fit by REML
## ***** sommer 3.3 *****
## =====
##           logLik      AIC      BIC Method Converge
## Value 167.0252 -322.0505 -298.5695  MNR      TRUE
## =====
## Variance-Covariance components:
##           VarComp VarCompSE Zratio
## Name.Yield-Yield      3.7090  1.68155  2.206
## Name.Yield-Weight      0.9071  0.37953  2.390
## Name.Weight-Weight      0.2243  0.08777  2.556
## Env:Name.Yield-Yield      5.0922  1.47906  3.443
## Env:Name.Yield-Weight      1.0269  0.30773  3.337
## Env:Name.Weight-Weight      0.2101  0.06662  3.153
## units.Yield-Yield      4.3837  0.64951  6.749
## units.Yield-Weight      0.9077  0.14147  6.416
## units.Weight-Weight      0.2280  0.03378  6.751
## =====
## Fixed effects:
##
```

```

## $Yield
##           Estimate Std. Error  t value
## (Intercept) 14.741988  0.6783190 21.733118
## EnvCA.2012  -3.199176  0.7474089 -4.280355
## EnvCA.2013  -4.003356  0.7850500 -5.099491
##
## $Weight
##           Estimate Std. Error  t value
## (Intercept)  0.5847384  0.1497086  3.905845
## EnvCA.2012  -0.9711518  0.1592560 -6.098053
## EnvCA.2013  -1.1643241  0.1681075 -6.926068
##
## =====
## Groups and observations:
##      Observ Groups
## Name      185     41
## Env:Name   185    123
## =====
## Use the '$' sign to access results and parameters

```

You may notice that we have added the `us(trait)` behind the random effects. This is to indicate the structure that should be assumed in the multivariate model. The `diag(trait)` used behind a random effect (i.e. Name) indicates that for the traits modeled (Yield and Weight) there's no a covariance component and should not be estimated, whereas `us(trait)` assumes that for such random effect, there's a covariance component to be estimated (i.e. covariance between Yield and Weight for the random effect Name). Same applies for the residual part (`rcov`).

5) Multivariate heterogeneous variance models

This is just an extension of the univariate heterogeneous variance models but at the multivariate level. This would be a CS+DIAG multivariate model:

```

data(example)
head(example)

```

```

##           Name      Env Loc Year      Block Yield      Weight
## 33  Manistee(MSL292-A) CA.2013  CA 2013  CA.2013.1      4 -1.904711
## 65           CO02024-9W 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 <- mmer2(cbind(Yield, Weight) ~ Env,
              random= ~ us(trait):Name + us(trait):at(Env):Name,
              rcov= ~ us(trait):at(Env):units,
              data=example, silent = TRUE)

```

```

summary(ans1)

```

```

## =====
##      Multivariate Linear Mixed Model fit by REML
## ***** sommer 3.3 *****
## =====
##           logLik      AIC      BIC Method Converge
## Value 177.8154 -343.6309 -320.1498  MNR      TRUE

```



```

## =====
## Variance-Covariance components:
##                               VarComp VarCompSE Zratio
## Name.Yield-Yield             3.32273   1.45378 2.2856
## Name.Yield-Weight            0.79471   0.32646 2.4343
## Name.Weight-Weight           0.19102   0.07508 2.5442
## CA.2011:Name.Yield-Yield     8.69977   4.01009 2.1695
## CA.2011:Name.Yield-Weight    1.77759   0.83838 2.1203
## CA.2011:Name.Weight-Weight   0.35940   0.17886 2.0094
## CA.2012:Name.Yield-Yield     2.57328   1.95113 1.3189
## CA.2012:Name.Yield-Weight    0.33269   0.39868 0.8345
## CA.2012:Name.Weight-Weight   0.03843   0.08601 0.4468
## CA.2013:Name.Yield-Yield     5.46662   2.16187 2.5287
## CA.2013:Name.Yield-Weight    1.34663   0.50455 2.6689
## CA.2013:Name.Weight-Weight   0.32893   0.12203 2.6954
## CA.2013:units.Yield-Yield    2.56131   0.63996 4.0023
## CA.2013:units.Yield-Weight   0.44569   0.12645 3.5246
## CA.2013:units.Weight-Weight  0.12232   0.03057 4.0009
## CA.2011:units.Yield-Yield    4.93845   1.52314 3.2423
## CA.2011:units.Yield-Weight   0.99446   0.32150 3.0932
## CA.2011:units.Weight-Weight  0.23982   0.07394 3.2433
## CA.2012:units.Yield-Yield    5.73843   1.31505 4.3637
## CA.2012:units.Yield-Weight   1.27999   0.30150 4.2454
## CA.2012:units.Weight-Weight  0.31804   0.07285 4.3657
## =====
## Fixed effects:
##
## $Yield
##           Estimate Std. Error  t value
## (Intercept) 14.498149  0.7889099 18.377447
## EnvCA.2012  -3.009529  0.8264130 -3.641677
## EnvCA.2013  -3.731619  0.8754603 -4.262465
##
## $Weight
##           Estimate Std. Error  t value
## (Intercept)  0.5746021  0.1682650  3.414865
## EnvCA.2012  -0.9334375  0.1697682 -5.498307
## EnvCA.2013  -1.1375573  0.1914174 -5.942811
##
## =====
## Groups and observations:
##           Observ Groups
## Name           185    41
## CA.2011:Name   185    41
## CA.2012:Name   185    41
## CA.2013:Name   185    41
## =====
## Use the '$' sign to access results and parameters

```

Any number of random effects can be specified with different structures.

6) Including special functions

Several random effects require the use of covariance structures that specify an special relationship among the levels of such random effect. The sommer package includes the g() function to include such known covariance structures:

```
data(example)
head(example)
```

```
##           Name      Env Loc Year      Block Yield      Weight
## 33  Manistee(MSL292-A) CA.2013  CA 2013  CA.2013.1      4 -1.904711
## 65           CO02024-9W 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
```

```
K[1:4,1:4]
```

```
##           Manistee(MSL292-A) CO02024-9W MSL007-B MSR169-8Y
## Manistee(MSL292-A)                1          0          0          0
## CO02024-9W                        0          1          0          0
## MSL007-B                          0          0          1          0
## MSR169-8Y                         0          0          0          1
```

```
ans1 <- mmer2(Yield ~ Env,
              random= ~ g(Name) + at(Env):g(Name),
              rcov= ~ at(Env):units,
              G=list(Name=K),
              data=example, silent = TRUE)
```

```
summary(ans1)
```

```
## =====
##           Multivariate Linear Mixed Model fit by REML
## ***** sommer 3.3 *****
## =====
##           logLik      AIC      BIC Method Converge
## Value -15.42982 36.85964 46.52071 MNR TRUE
## =====
## Variance-Covariance components:
##           VarComp VarCompSE Zratio
## g(Name).Yield-Yield      2.962  1.4964  1.980
## CA.2011:g(Name).Yield-Yield 10.148  4.5107  2.250
## CA.2012:g(Name).Yield-Yield  1.879  1.8700  1.005
## CA.2013:g(Name).Yield-Yield  6.629  2.5027  2.649
## CA.2013:units.Yield-Yield   2.560  0.6398  4.001
## CA.2011:units.Yield-Yield   4.942  1.5246  3.242
## CA.2012:units.Yield-Yield   5.725  1.3119  4.364
## =====
## Fixed effects:
##
## $Yield
##           Estimate Std. Error  t value
## (Intercept) 16.507676  0.8268629 19.964224
## EnvCA.2012  -5.816887  0.8575779 -6.782926
## EnvCA.2013  -6.412430  0.9356441 -6.853493
```

```
##
## =====
## Groups and observations:
##           Observ Groups
## g(Name)      185      41
## CA.2011:g(Name) 185      41
## CA.2012:g(Name) 185      41
## CA.2013:g(Name) 185      41
## =====
## Use the '$' sign to access results and parameters
```

and for multivariate models:

```
data(example)
head(example)
```

```
##           Name      Env Loc Year      Block Yield      Weight
## 33  Manistee(MSL292-A) CA.2013  CA 2013  CA.2013.1      4 -1.904711
## 65           CO02024-9W 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
```

```
K[1:4,1:4]
```

```
##           Manistee(MSL292-A) CO02024-9W MSL007-B MSR169-8Y
## Manistee(MSL292-A)           1           0           0           0
## CO02024-9W                   0           1           0           0
## MSL007-B                       0           0           1           0
## MSR169-8Y                       0           0           0           1
```

```
ans1 <- mmer2(cbind(Yield, Weight) ~ Env,
              random= ~ us(trait):g(Name) + us(trait):at(Env):g(Name),
              rcov= ~ us(trait):at(Env):units,
              G=list(Name=K),
              data=example, silent = TRUE)
summary(ans1)
```

```
## =====
##           Multivariate Linear Mixed Model fit by REML
## ***** sommer 3.3 *****
## =====
##           logLik      AIC      BIC Method Converge
## Value 177.8154 -343.6309 -320.1498   MNR      TRUE
## =====
## Variance-Covariance components:
##           VarComp VarCompSE Zratio
## g(Name).Yield-Yield      3.32273  1.45378 2.2856
## g(Name).Yield-Weight      0.79471  0.32646 2.4343
## g(Name).Weight-Weight      0.19102  0.07508 2.5442
## CA.2011:g(Name).Yield-Yield 8.69977  4.01009 2.1695
## CA.2011:g(Name).Yield-Weight 1.77759  0.83838 2.1203
## CA.2011:g(Name).Weight-Weight 0.35940  0.17886 2.0094
## CA.2012:g(Name).Yield-Yield 2.57328  1.95113 1.3189
## CA.2012:g(Name).Yield-Weight 0.33269  0.39868 0.8345
## CA.2012:g(Name).Weight-Weight 0.03843  0.08601 0.4468
```

```

## CA.2013:g(Name).Yield-Yield  5.46662  2.16187 2.5287
## CA.2013:g(Name).Yield-Weight  1.34663  0.50455 2.6689
## CA.2013:g(Name).Weight-Weight 0.32893  0.12203 2.6954
## CA.2013:units.Yield-Yield    2.56131  0.63996 4.0023
## CA.2013:units.Yield-Weight   0.44569  0.12645 3.5246
## CA.2013:units.Weight-Weight  0.12232  0.03057 4.0009
## CA.2011:units.Yield-Yield    4.93845  1.52314 3.2423
## CA.2011:units.Yield-Weight   0.99446  0.32150 3.0932
## CA.2011:units.Weight-Weight  0.23982  0.07394 3.2433
## CA.2012:units.Yield-Yield    5.73843  1.31505 4.3637
## CA.2012:units.Yield-Weight   1.27999  0.30150 4.2454
## CA.2012:units.Weight-Weight  0.31804  0.07285 4.3657
## =====
## Fixed effects:
##
## $Yield
##           Estimate Std. Error  t value
## (Intercept) 14.498149  0.7889099 18.377447
## EnvCA.2012  -3.009529  0.8264130 -3.641677
## EnvCA.2013  -3.731619  0.8754603 -4.262465
##
## $Weight
##           Estimate Std. Error  t value
## (Intercept)  0.5746021  0.1682650  3.414865
## EnvCA.2012  -0.9334375  0.1697682 -5.498307
## EnvCA.2013  -1.1375573  0.1914174 -5.942811
##
## =====
## Groups and observations:
##           Observ Groups
## g(Name)      185      41
## CA.2011:g(Name) 185      41
## CA.2012:g(Name) 185      41
## CA.2013:g(Name) 185      41
## =====
## Use the '$' sign to access results and parameters

```

Notice that the `g()` function is applied at the random effect called “Name”, and the covariance structure is provided in the argument “G”. In the example, we used a diagonal covariance structure for demonstration purposes but any dense covariance matrix can be used.

Other special functions such as `overlay()` for overlay models, `eig()` for an eigen decomposition of the covariance matrix, `grp()` for customized random effects providing an incidence matrix, and `sp12D()` are available. Take a look at the help page for each of these special functions.

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

We will use the CPdata to show the use of 2-dimensional splines for accommodating 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 accommodate 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(CPdata)
head(CPpheno)

##      id Row Col Year      color  Yield FruitAver Firmness Rowf Colf
## P003 P003  3  1 2014 0.10075269 154.67   41.93  588.917   3   1
## P004 P004  4  1 2014 0.13891940 186.77   58.79  640.031   4   1
## P005 P005  5  1 2014 0.08681502  80.21   48.16  671.523   5   1
## P006 P006  6  1 2014 0.13408561 202.96   48.24  687.172   6   1
## P007 P007  7  1 2014 0.13519278 174.74   45.83  601.322   7   1
## P008 P008  8  1 2014 0.17406685 194.16   44.63  656.379   8   1

CPgeno[1:4,1:4]

##      scaffold_50439_2381 scaffold_39344_153 uneak_3436043 uneak_2632033
## P003                    0                    0                0                1
## P004                    0                    0                0                1
## P005                    0                   -1                0                1
## P006                   -1                   -1               -1                0

#### create the variance-covariance matrix
A <- A.mat(CPgeno) # additive relationship matrix
#### look at the data and fit the model
head(CPpheno)

##      id Row Col Year      color  Yield FruitAver Firmness Rowf Colf
## P003 P003  3  1 2014 0.10075269 154.67   41.93  588.917   3   1
## P004 P004  4  1 2014 0.13891940 186.77   58.79  640.031   4   1
## P005 P005  5  1 2014 0.08681502  80.21   48.16  671.523   5   1
## P006 P006  6  1 2014 0.13408561 202.96   48.24  687.172   6   1
## P007 P007  7  1 2014 0.13519278 174.74   45.83  601.322   7   1
## P008 P008  8  1 2014 0.17406685 194.16   44.63  656.379   8   1

mix1 <- mmer2(Yield~1,
              random=~g(id)
                + Rowf + Colf
                + spl2D(Row,Col),
              rcov=~units,
              G=list(id=A), silent=TRUE,
              data=CPpheno)
summary(mix1)

## =====
##      Multivariate Linear Mixed Model fit by REML
## ***** sommer 3.3 *****
## =====
##      logLik      AIC      BIC Method Converge
## Value -151.201 304.4021 308.2937   MNR      TRUE
## =====
## Variance-Covariance components:
##      VarComp VarCompSE Zratio
## g(id).Yield-Yield      782.9    318.9 2.4549
```

```

## Rowf.Yield-Yield           814.8    390.9 2.0846
## Colf.Yield-Yield           182.1    129.6 1.4056
## at_FIELD1_2Dspl.Yield-Yield 514.0    694.8 0.7397
## units.Yield-Yield          2922.8    294.2 9.9360
## =====
## Fixed effects:
##
## $Yield
##      Estimate Std. Error  t value
## Intercept 132.1423   8.791225 15.03116
##
## =====
## Groups and observations:
##      Observ Groups
## g(id)      362    363
## Rowf       362    13
## Colf       362    36
## at_FIELD1_2Dspl 362    168
## =====
## Use the '$' sign to access results and parameters

```

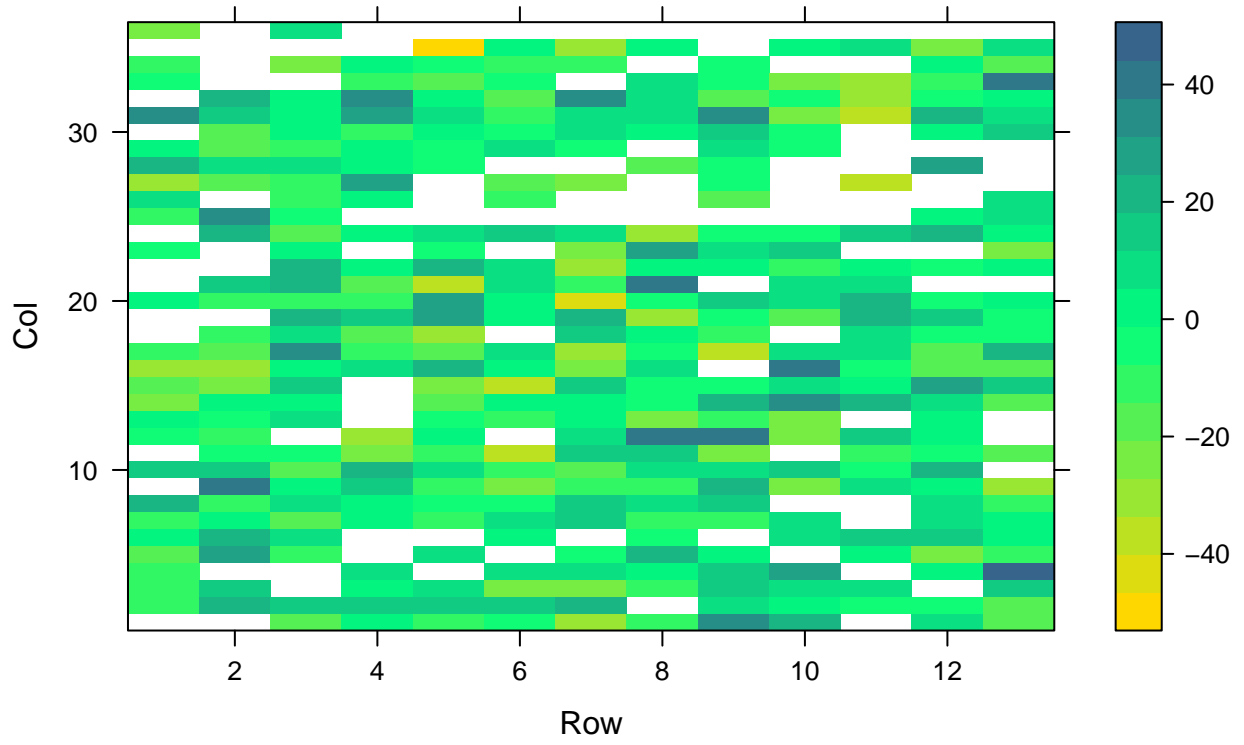
Notice that the job is done by the `sp12D()` function that takes the Row and Col information to fit a spatial kernel. When multiple fields are available the function has an additional argument call `at` which allows to fit a different spatial kernel in each field. For example if there was multiple fields the use of the random call would look like:

```
random=~ sp12D(Row,Col, at=FIELD)
```

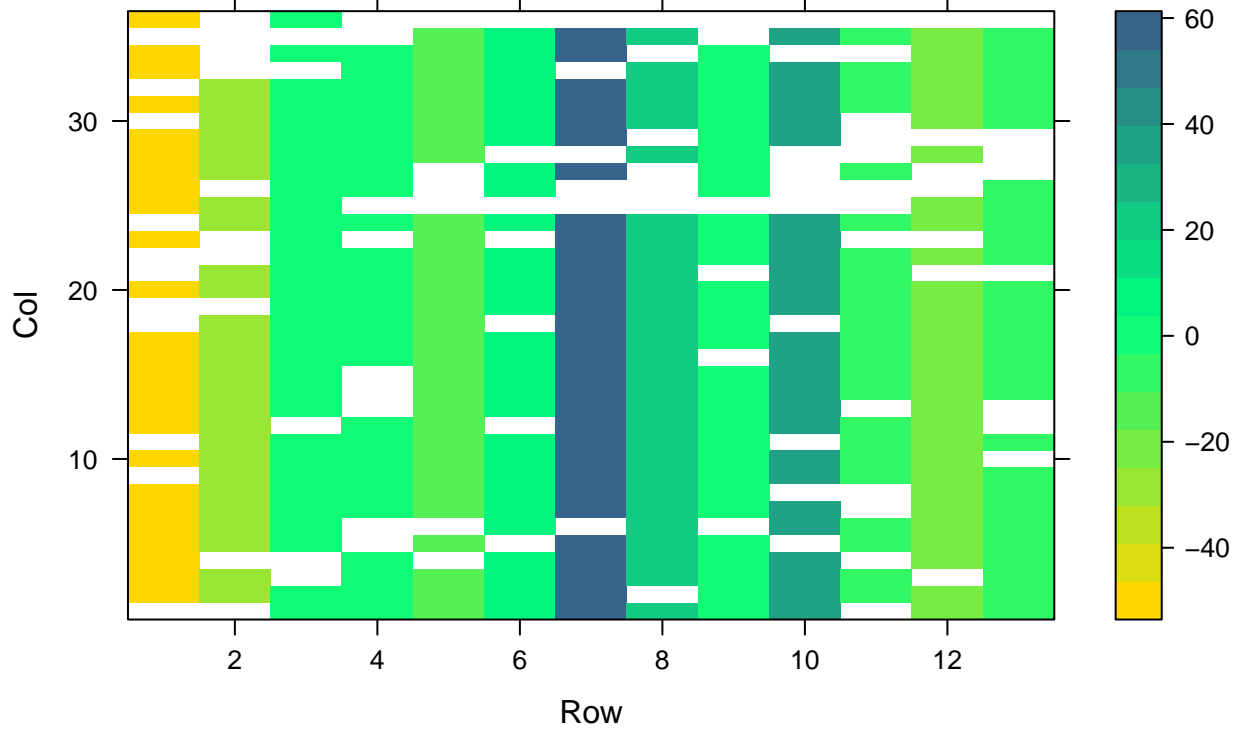
where FIELD would correspond to the name of the column of the dataset where the identifier for the different environments is.

```
#### get the spatial plots
fittedvals <- spatPlots(mix1,row = "Row", range = "Col")
```

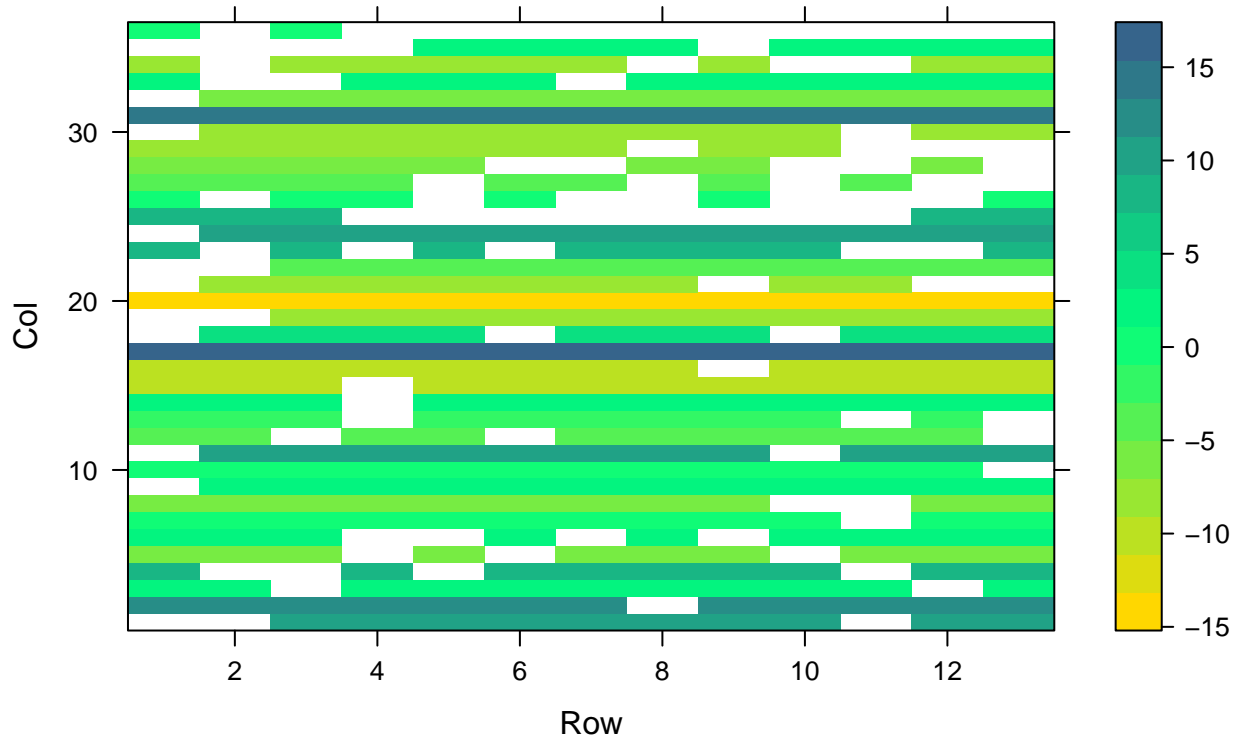
fit_g.id.



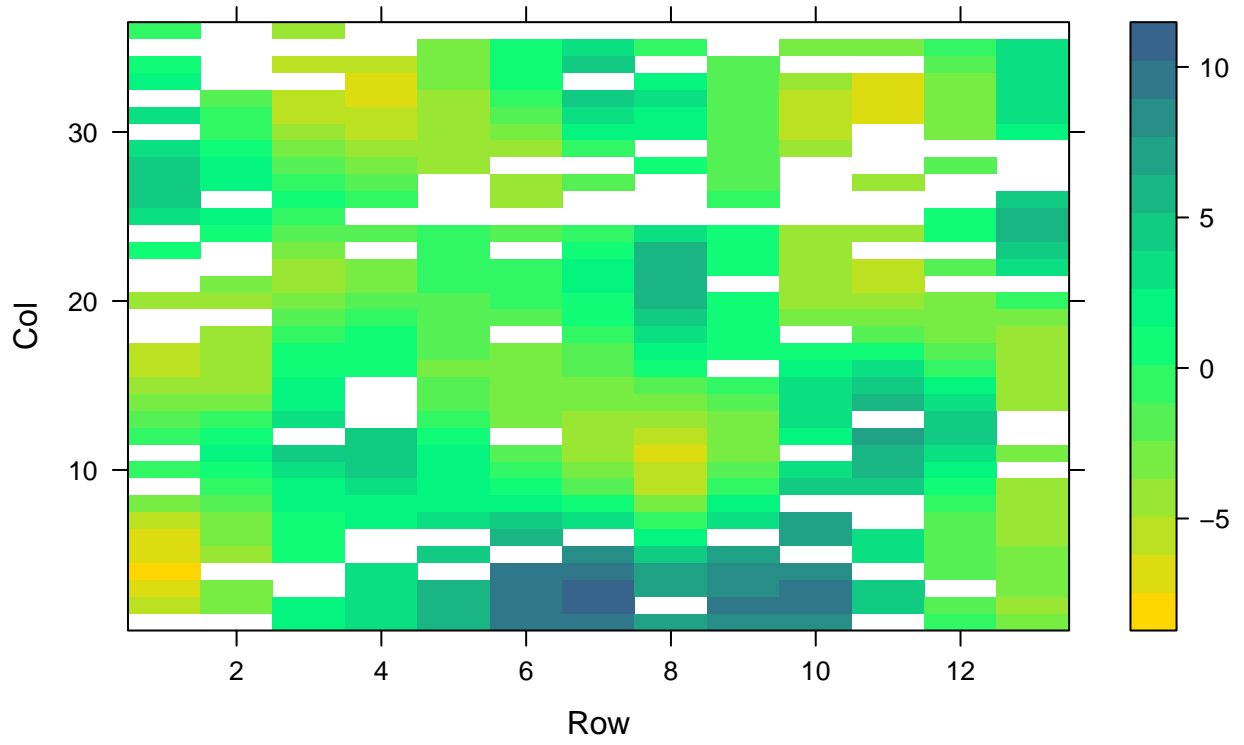
fit_Rowf

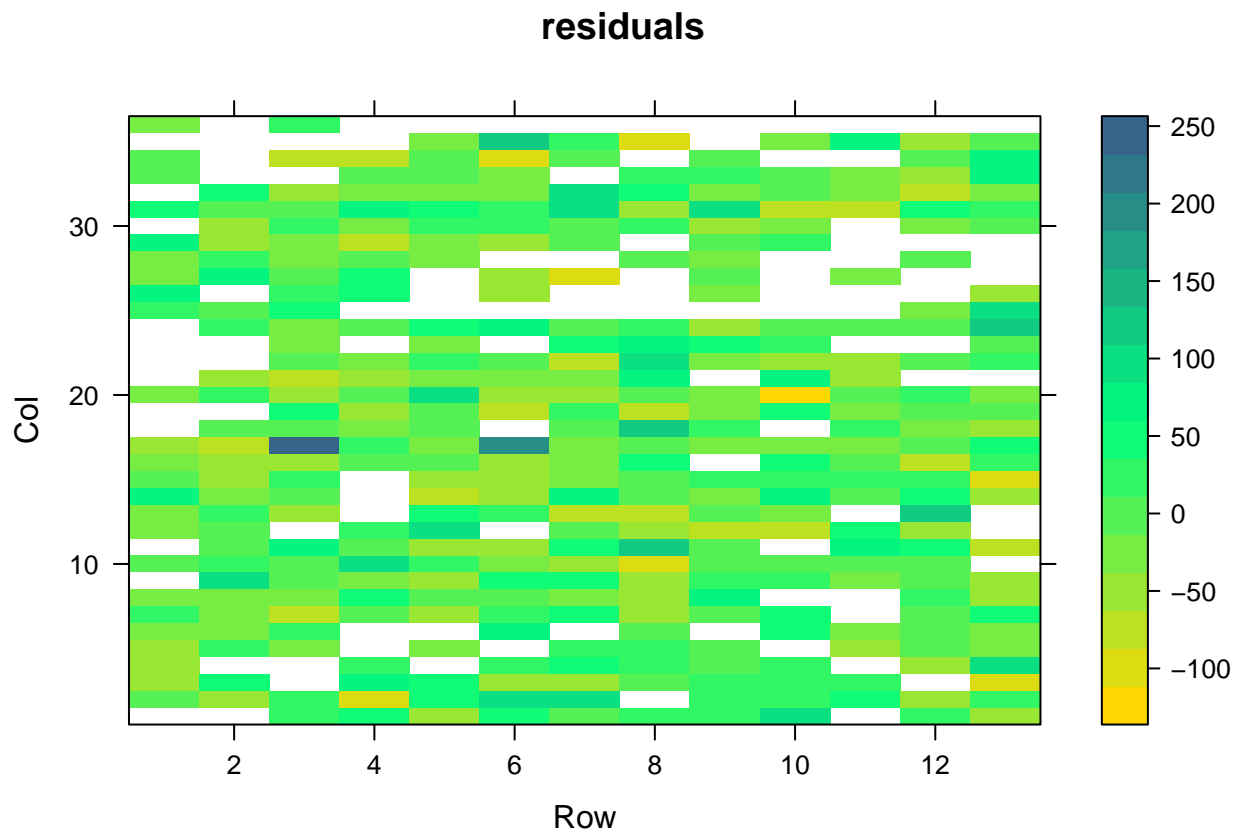


fit_Colf



fit_spl2D





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.
- Rodríguez-Álvarez, María Xosé, 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.
- Tunncliffe W. 1989. On the use of marginal likelihood in time series model estimation. *JRSS* 51(1):15-27.