# Multitrait models in BGLR

Paulino Pérez-Rodríguez [1]
Gustavo de los Campos [2]

[1]Colegio de Posgraduados [2]Michigan State University
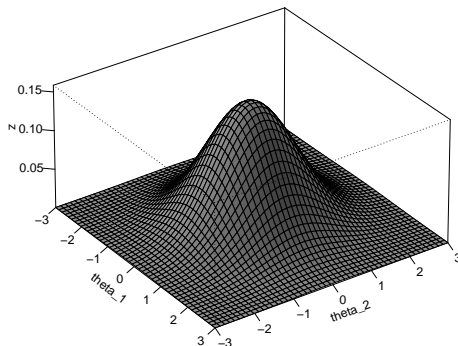
June, 2020.

# Contents

# Multivariate Normal distribution

The MN distribution is a generalization of the univariate normal distribution to higher dimensions. If $\boldsymbol{\theta} \sim MN(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, where $\boldsymbol{\mu}$ is the mean and $\boldsymbol{\Sigma}$ is the variance covariance matrix. The density function of $\boldsymbol{\theta}$ is given by:

$$p(\boldsymbol{\theta}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{(2\pi)^{k/2}|\boldsymbol{\Sigma}|^{1/2}} \exp\left\{-\frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\mu})'\boldsymbol{\Sigma}(\boldsymbol{\theta} - \boldsymbol{\mu})\right\}$$

# Matrix normal distribution or matrix Gaussian distribution

The MG distribution is a generalization of the MN distribution. The random matrix $\Theta_{n \times q}$ follows a MG distribution $MG(\boldsymbol{M}, \boldsymbol{V}, \boldsymbol{W})$ and the density function is given by:

$$p(\Theta|\boldsymbol{M}, \boldsymbol{V}, \boldsymbol{W}) = \frac{\exp\{-0.5 tr[\boldsymbol{W}^{-1}(\Theta - \boldsymbol{M})'\boldsymbol{V}^{-1}(\Theta - \boldsymbol{M})]\}}{(2\pi)^{nq/2}|\boldsymbol{W}|^{n/2}|\boldsymbol{V}|^{q/2}},$$

where $\boldsymbol{M}$ is a location parameter, a $n \times q$ matrix, $\boldsymbol{V}$ is a $n \times n$ matrix, a scale parameter for rows and $\boldsymbol{W}$ is a $q \times q$, a scale parameter for columns.

It can be shown that $\Theta_{n \times t} \sim MG(\boldsymbol{M}, \boldsymbol{V}, \boldsymbol{W})$, if and only if $vec(\Theta_{n \times q}) \sim MN(vec(\boldsymbol{M}), \boldsymbol{W} \otimes \boldsymbol{V})$.

# Inverse Wishart distribution

$\Theta \sim IW_q(\boldsymbol{S}, df)$, where $\boldsymbol{S}_{q \times q}$ is a scale matrix, which should be positive definite, $df$ corresponds to the degrees of freedom, $df > q - 1$. The density function is proportional to:

$$p(\Theta | \boldsymbol{S}, df) \propto |\boldsymbol{S}|^{df/2} |\Theta|^{-(df+q+1)/2} \exp\{-0.5 tr(\boldsymbol{S}\Theta^{-1})\}$$

The expectation is defined if $df > q + 1$.

# Multivariate Multiple Linear Regression

Let $\boldsymbol{y}_i = (y_{i1}, ..., y_{iq})'$, $i = 1, ..., n$ (individuals) and $k = 1, ..., q$ (traits), and let $\boldsymbol{X}$ represent a matrix of covariates of dimension $n \times p$ and $\boldsymbol{B} = (\boldsymbol{\beta}_1, ..., \boldsymbol{\beta}_q)$ a matrix of regression coefficients. Without loss of generality, assuming that the the phenotypes have been centered by trait, the MMLR model can be written as:

$$\begin{pmatrix} y_{11} & \cdots & y_{1q} \\ \vdots & \ddots & \vdots \\ y_{n1} & \cdots & y_{nq} \end{pmatrix} = \begin{pmatrix} x_{11} & \cdots & x_{1p} \\ \vdots & \ddots & \vdots \\ x_{n1} & \cdots & x_{np} \end{pmatrix} \begin{pmatrix} \beta_{11} & \cdots & \beta_{1q} \\ \vdots & \ddots & \vdots \\ \beta_{p1} & \cdots & \beta_{pq} \end{pmatrix} + \begin{pmatrix} e_{11} & \cdots & e_{1q} \\ \vdots & \ddots & \vdots \\ e_{n1} & \cdots & e_{nq} \end{pmatrix}$$

or in compact matrix notation:

$$\boldsymbol{Y} = \boldsymbol{X}\boldsymbol{B} + \boldsymbol{E}$$

here for example $\boldsymbol{E} \sim MG(\boldsymbol{0}, \boldsymbol{I}, \boldsymbol{R}_0)$, with $\boldsymbol{R}_0$ matrix of dimensions $q \times q$.

# Multivariate Linear Mixed Model (Bayesian context)

Without loss of generality consider:

$$\boldsymbol{Y} = \boldsymbol{1}\boldsymbol{\mu}' + \boldsymbol{X}\boldsymbol{B} + \boldsymbol{U} + \boldsymbol{E},$$

with $\boldsymbol{Y}$, $\boldsymbol{X}$, $\boldsymbol{B}$, $\boldsymbol{E}$ as defined before, $\boldsymbol{\mu} = (\mu_1, ..., \mu_q)'$ a vector of intercepts for traits, and

$$\boldsymbol{U} = \begin{pmatrix} u_{11} & \cdots & u_{1t} \\ \vdots & \ddots & \vdots \\ u_{n1} & \cdots & u_{nt} \end{pmatrix}$$

is a matrix of random effects.

# Continue...

With:

1. $E \sim MG(\mathbf{0}, I, R_0)$, $R_0$ (unknown) covariance matrix for residuals of dim $q \times q$.
2. $p(\mu) \propto 1$.
3. $B \sim MG(\mathbf{0}, I, \Omega)$, with $\Omega$ (unknown) covariance matrix for regression coefficients.
4. $U \sim MG(\mathbf{0}, K, G)$, $K$ (known) covariance matrix to model relationship between individuals and $G$ (unknown) covariance matrix to model within subject relations.

We assign Inverse Wishart distributions to $R_0$, $\Omega$ and $G$.

# Continue...

Posterior distribution has no closed form, but all full conditionals necessary to implement a Gibbs sampler (Geman and Geman, 1984) have closed form. We need to sample from:

1. $\mu|else$.
2. $B|else$.
3. $U|else$.
4. $R_0|else$.
5. $\Omega|else$.
6. $G|else$.

# Continue...

We consider several modeling strategies for $R_0$, $\Omega$ and $G$. Without loss of generality consider the case of $\Omega$.

| Structure | Prior |
|---|---|
| | Hyper-parameters |
| Unstructured | $\Omega \sim IW(\Omega | S_0, df_0)$; |
| | $S_0$: prior scale matrix ($q \times q$) |
| | $df_0$: prior degree of freedom |
| Diagonal | $\Omega = \mathrm{diag}(\omega_{11}, \ldots, \omega_{tt})$, |
| | $p(\boldsymbol{\omega}) = \prod_{j=1}^{q} \chi^{-2}(\omega_{jj} | S_{0jj}, df_{0j})$ |
| | $S_0$: prior scale vector ($q \times 1$) |
| | $df_0$: prior $df$ vector ($q \times 1$) |
| Factor analytic | $\Omega = \Lambda\Lambda' + \Psi$; $\Lambda$:factor loadings, $\Psi = \mathrm{diag}(\psi_{jj})$ |
| | $p(\Lambda, \Psi) = MVN(vec(\Lambda) | \mathbf{0}, \sigma^2 \mathbf{I}) \times \prod_{j=1}^{q} \chi^{-2}(\psi_{jj} | S_{0jj}, df_{0j})$ |
| | $S_0$: prior scale vector ($q \times 1$) |
| | $df_0$: prior $df$ vector ($q \times 1$); $\sigma^2$: prior variance. |
| Recursive | $\Omega = (\mathbf{I} - \mathbf{B})^{-1}\Psi((\mathbf{I} - \mathbf{B})^{-1})'$, $B_{q \times q}$ regression coefficients; $\Psi = \mathrm{diag}(\psi_{jj})$. |
| | $p(\mathbf{B}, \Psi) = MVN(vec(\mathbf{B}) | \mathbf{0}, \sigma^2 \mathbf{I}) \times \prod_{j=1}^{q} \chi^{-2}(\psi_{jj} | S_{0jj}, df_{0j})$ |
| | $S_0$: prior scale vector ($q \times 1$) |
| | $df_0$: prior $df$ vector ($q \times 1$); $\sigma^2$: prior variance. |

# Continue...

We can also assign different prior distribution to components of **B**, for example mixture of normals (e.g. Cheng et al., 2018) or flat priors.



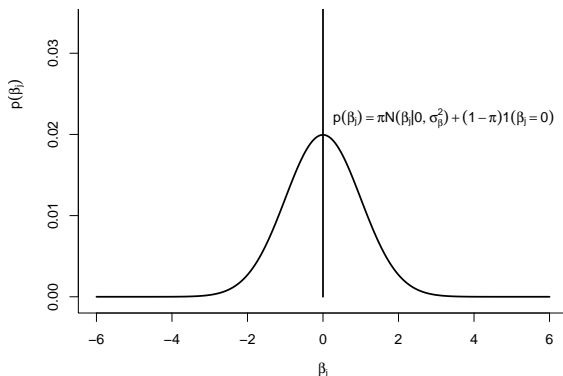$$p(\beta_j) = \pi N(\beta_j|0, \sigma_\beta^2) + (1-\pi)1(\beta_j = 0)$$

Figure 1: Spike-Slab.

# Multitrait BGLR

BGLR now includes the function Multitrait that fits a MLMM. The function is a re-implementation/extension of the function MTM included in the MTM package (http://quantgen.github.io/MTM/vignette.html).

Key features:

1. Heavy computational routines implemented in the C programming language.
2. Gaussian and Gaussian mixtures priors for regression coefficients.
3. Different covariance structures for $R_0$, $\Omega$ and $G$ (UNstructured, DIAGonal, Factor Analytic, RECursive).
4. Arbitrary number of random effects.

Visit https://github.com/gdlc/BGLR-R.

## Continue...

The Multitrait function fits the model:

$$Y = 1\mu' + X_1B_1 + X_2B_2 + \cdots + U_1 + U_2 + \cdots + E$$

where all the terms have been defined before. The function implements a Gibbs sampler with scalar update. Missing values are imputed at each iteration of Gibbs algorithm.

We will discuss only the more basic models, but we present the information required to fit more complex models, the sample code is accesible in the github website.

# Model specification

The user interface is similar to the BGLR function and is list based.

### Table 2. Specifying components of the linear predictor.

| Component | Specification |
|---|---|
| Predictors with regression coefficients with flat prior | `list(X=?,model="FIXED")` |
| Predictors with regression coefficients with Gaussian prior | `list(X=?,model="BRR",Cov=?)` |
| Random effects | `list(K=?,model="RKHS",Cov=?)` |
| Predictors with regression coefficients with Spike Slab prior | `list(X=?,model="SSVS",Cov=?, inclusionProb=?)` |
| | |
| | `inclusionProb` is a list with two components, `probIn`: a vector with inclusion probabilities and `counts`: a vector of prior counts. |

# Continue...

## Table 3. Specifying covariance structures in R code.

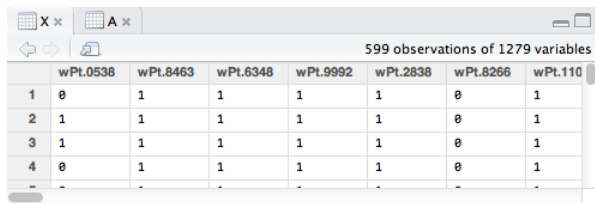| Covariance structure | Specification |
| --- | --- |
| Unstructured | `Cov=list(type="UN",df0=?,S0=?)` |
| Diagonal | `Cov=list(type="DIAG",df0=?,S0=?)` |
| Factor analytic | `Cov=list(type="FA",df0=?,S0=?, M=?, var=?)` |
| | `M` is a logical matrix with $t$ rows and `TRUE` for loadings that the user want to estimate and `FALSE` for those that should be zeroed out. |
| Recursive | `Cov=list(type="FA",df0=,S0=?, M=?, var=?)` |
| | `M` is a logical matrix of dimensions $t \times t$, diagonal must be set to `FALSE` and lower diagonal elements are set to `TRUE` to model recursion between traits. |

# Example 1: Fitting RR

Data for $n = 599$ wheat lines evaluated in 4 environments, wheat improvement program, CIMMyT. The dataset includes $p = 1279$ molecular markers ($x_{ij}, i = 1, ..., n, \ j = 1, ..., p$) (coded as 0,1). The pedigree information is also available.
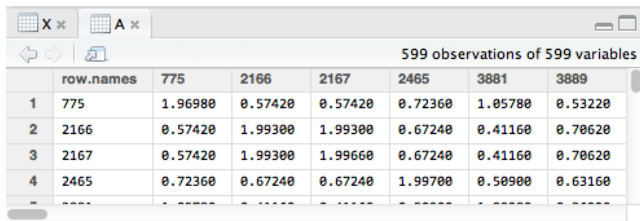
Lets load the dataset in R,

1. Load R.
2. Install BGLR package (if not yet installed).
3. Load the package.
4. Load the data.

# Continue...

You can explore the MM matrix, pedigree matrix within R,

We fit the model:

$$Y = \mathbf{1}\mu' + XB + E,$$

$B \sim MG(\mathbf{0}, I, \Omega)$, with $\Omega$ (unknown) covariance matrix for effects. We assign an inverse Whishart distribution to $R_0$ and $\Omega$ (default settings in the software).

```
#Loading software
library(BGLR)

#Loading data
data(wheat)
y<-wheat.Y
X<-wheat.X
X<-scale(X)/sqrt(ncol(X))

#Linear predictor
ETA<-list(list(X=X,model="BRR"))

#Model fitting
set.seed(123)
fm<-Multitrait(y=y,ETA=ETA,nIter=10000,burnIn=5000)
```

# Continue...

```
#Residual covariance matrix
fm$resCov

#Genetic covariance matrix
fm$ETA[[1]]$Cov

#Marker effects
fm$ETA[[1]]$beta

#Edited output
>fm$resCov
...
$type
[1] "UN"

$R
          1          2          4          5
1  0.53315070 0.07797869 -0.1055548 0.05279538
2  0.07797869 0.57551110  0.3034339 0.16497247
4 -0.10555484 0.30343389  0.6196055 0.11723637
5  0.05279538 0.16497247  0.1172364 0.59673421
```

# Continue...

```
#Edited output
> fm$ETA[[1]]$Cov
...
$type
[1] "UN"

$Omega
            1          2          4          5
1  0.5632461 -0.1146953 -0.1213392 -0.2037660
2 -0.1146953  0.5058328  0.4081491  0.2661113
4 -0.1213392  0.4081491  0.4966503  0.2942788
5 -0.2037660  0.2661113  0.2942788  0.4959512
...
```

## Equivalent code to specify covariance structures:

```
ETA<-list(list(X=X,model="BRR",Cov=list(type="UN")))
residual<-list(type="UN")

#Model fitting
set.seed(123)
fm<-Multitrait(y=y,ETA=ETA,resCov=residual,nIter=10000,burnIn=5000)
```

# Example 2: Random effect model with unstructured covariance matrices

We fit the model:

$$Y = 1\mu' + U + E,$$

$U \sim MG(0, K, G)$, $K$ (known) covariance matrix to model relationship between individuals and $G$ (unknown) covariance matrix to model within subject relations. We assign an inverse Whishart distribution to $R_0$ and $G$ (default settings in the software).

```
library(BGLR)
data(wheat)
K<-wheat.A
y<-wheat.Y

ETA<-list(list(K=K,model="RKHS"))

#Fit model
set.seed(123)
fm<-Multitrait(y=y,ETA=ETA,nIter=10000,burnIn=5000)
```

## Continue...

```
#Residual covariance matrix
fm$resCov

#Genetic covariance matrix
fm$ETA[[1]]$Cov

#Random effects
fm$ETA[[1]]$u

#Edited output
>fm$resCov
...
$type
[1] "UN"
...
$R
            1          2          4          5
1  0.567446646 0.0343702 -0.1180475 0.004740618
2  0.034370198 0.5645890  0.2450678 0.151970705
4 -0.118047470 0.2450678  0.5043002 0.103934039
5  0.004740618 0.1519707  0.1039340 0.531484892
```

# Continue...

```
#Edited output
>fm$ETA[[1]]$Cov
...
$type
[1] "UN"
...
$G
            1           2           4           5
1  0.30324330 -0.04113035 -0.06129195 -0.08841091
2 -0.04113035  0.28808179  0.26584470  0.15342103
4 -0.06129195  0.26584470  0.34676292  0.17047220
5 -0.08841091  0.15342103  0.17047220  0.30945972
```

# Questions?

# References

📄 Geman, S. and Geman, D. (1984).
Stochastic relaxation, Gibbs distributions, and the Bayesian restoration of images.
*IEEE Transactions on pattern analysis and machine intelligence*, 6, 6, 721–741.

📕 Sorensen, D. and Gianola, D. (2007).
*Likelihood, Bayesian, and MCMC methods in quantitative genetics*.
Springer Science & Business Media.