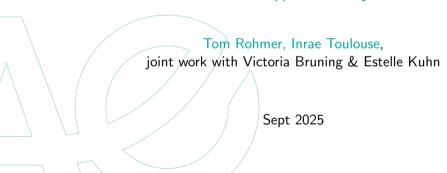


copula model to improve genetic parameter estimation for dependent traits.

StatMathAppli, 2025, Fréjus



> multitrait G+E animal model

⊳ Every phenotypic observation on an animal is determined by environmental and genetic factors and may be defined by the following model:

Phenotypic observation

= envir. effects + genetic effects + resid. effects





bivariate mixed model (G+E)

With bivariate phenotypes, the most classical inference model in genetic is

$$\begin{cases} \mathbf{y}_1 = X_1 \boldsymbol{\beta}_1 + Z_1 \mathbf{a}_1 + \varepsilon_1 \\ \mathbf{y}_2 = X_2 \boldsymbol{\beta}_2 + Z_2 \mathbf{a}_2 + \varepsilon_2. \end{cases}$$

- y_i the phenotype vectors (neither identically distributed nor independent)
- β_i parameter vectors to estimate, X_i design matrices related to fixed effects,
- a_i unobservable (latent) vectors, Z_i incidence matrices related to genetic effects
- ε_i residual vectors with components assumed i.i.d.





bivariate mixed model (G+E)

With bivariate phenotypes, the most classical inference model in genetic is

$$\begin{cases} y_1 = X_1 \beta_1 + Z_1 a_1 + \varepsilon_1 \\ y_2 = X_2 \beta_2 + Z_2 a_2 + \varepsilon_2. \end{cases}$$

 y_j the phenotype vectors (neither identically distributed nor independent) β_j parameter vectors to estimate, X_j design matrices related to fixed effects, a_j unobservable (latent) vectors, Z_j incidence matrices related to genetic effects ε_j residual vectors with components assumed i.i.d.

Particularly, the genetic effets (referred to as breeding values) are

$$a_{i,j} = 0.5(a_{i_S,j} + a_{i_D,j}) + M_{i,j},$$

where $a_{i_S,j}$ and $a_{i_D,j}$ are the BVs of the sire and dam $M_{i,j}$ are the Mendelian sampling terms with distribution





From the classical mixed model to Copula mixed model

The variance of the random terms (latent) is

$$var(\mathbf{a}_1, \mathbf{a}_2) = G \otimes A = \begin{pmatrix} \sigma_{\mathbf{a}_1}^2 A & \sigma_{\mathbf{a}_{12}} A \\ \sigma_{\mathbf{a}_{12}} A & \sigma_{\mathbf{a}_{2}}^2 A \end{pmatrix}, \quad A \text{ the kinship matrix,}$$





From the classical mixed model to Copula mixed model

The variance of the random terms (latent) is

$$var(\mathbf{a}_1, \mathbf{a}_2) = G \otimes A = \begin{pmatrix} \sigma_{\mathbf{a}_1}^2 A & \sigma_{\mathbf{a}_{12}} A \\ \sigma_{\mathbf{a}_{12}} A & \sigma_{\mathbf{a}_{2}}^2 A \end{pmatrix}, \quad A \text{ the kinship matrix,}$$

and for $i = 1, \ldots, n$

$$arepsilon_{i,j} \sim \mathcal{N}(0, \sigma_j^2), \quad (arepsilon_{i,1}, arepsilon_{i,2})$$
 is assumed jointly Gaussian

Then, REML precedures are used to estimate parameters.





From the classical mixed model to Copula mixed model

The variance of the random terms (latent) is

$$var(\mathbf{a}_1, \mathbf{a}_2) = G \otimes A = \begin{pmatrix} \sigma_{\mathbf{a}_1}^2 A & \sigma_{\mathbf{a}_{12}} A \\ \sigma_{\mathbf{a}_{12}} A & \sigma_{\mathbf{a}_{2}}^2 A \end{pmatrix}, \quad A \text{ the kinship matrix,}$$

and for $i = 1, \ldots, n$

$$arepsilon_{i,j} \sim \mathcal{N}(0, \sigma_j^2), \quad (arepsilon_{i,1}, arepsilon_{i,2})$$
 is assumed jointly Gaussian

Then, REML precedures are used to estimate parameters.

→ The multivariate-Gaussian assumption is often unrealistic due to the dependence structure of the observations, we propose the generalization:

$$arepsilon_{i,j} \sim \mathcal{N}(0, \sigma_j^2), \quad (arepsilon_{i,1}, arepsilon_{i,2}) \quad ext{has copula } C_{ heta}$$



Copulas

Definition: A copula $C: [0,1]^d \to [0,1]$ is the multivariate cumulative distribution function (c.d.f.) of a random vector whose marginal distributions are uniforms on [0,1].





Definition: A copula $C: [0,1]^d \to [0,1]$ is the multivariate cumulative distribution function (c.d.f.) of a random vector whose marginal distributions are uniforms on [0,1].

Theorem of [Sklar(1959)]

Let $\mathbf{X} = (X_1, \dots, X_d)$ be a d-dimensional random vector with c.d.f. \mathbf{F} and let F_1, \dots, F_d be the marginal c.d.f. of \mathbf{X} assumed continuous. Then it exists a unique copula C such that:

$$F(x) = C\{F_1(x_1), \dots, F_d(x_d)\}, \qquad x = (x_1, \dots, x_d) \in \mathbb{R}^d.$$





Definition: A copula $C: [0,1]^d \to [0,1]$ is the multivariate cumulative distribution function (c.d.f.) of a random vector whose marginal distributions are uniforms on [0,1].

Theorem of [Sklar(1959)]

Let $\mathbf{X} = (X_1, \dots, X_d)$ be a d-dimensional random vector with c.d.f. \mathbf{F} and let F_1, \dots, F_d be the marginal c.d.f. of \mathbf{X} assumed continuous. Then it exists a unique copula C such that:

$$F(x) = C\{F_1(x_1), \dots, F_d(x_d)\}, \qquad x = (x_1, \dots, x_d) \in \mathbb{R}^d.$$

 \triangleright The copula C characterizes the dependence structure of vector X.





Definition: A copula $C: [0,1]^d \to [0,1]$ is the multivariate cumulative distribution function (c.d.f.) of a random vector whose marginal distributions are uniforms on [0,1].

Theorem of [Sklar(1959)]

Let $\boldsymbol{X}=(X_1,\ldots,X_d)$ be a d-dimensional random vector with c.d.f. \boldsymbol{F} and let F_1,\ldots,F_d be the marginal c.d.f. of \boldsymbol{X} assumed continuous. Then it exists a unique copula C such that:

$$F(x) = C\{F_1(x_1), \dots, F_d(x_d)\}, \qquad x = (x_1, \dots, x_d) \in \mathbb{R}^d.$$

The copula C characterizes the dependence structure of vector X.



A. Sklar.

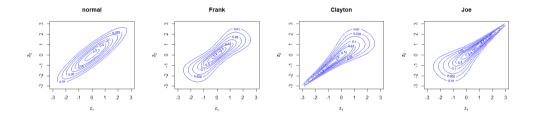
Fonctions de répartition à n dimensions et leurs marges.

Publications de l'Institut de Statistique de l'Université de Paris, 8:229-231, 1959.



INRA

Contour plots of bivariate distributions with Gaussian margins and several copula



The lack of consideration for an appropriate dependence structure (e.g., wrongly assuming a Gaussian distribution) may lead to poor estimation of variance parameters.



Rohmer, T., Ricard, A., David, I.

Copula miss-specification in REML multivariate genetic animal model estimation, Genetics Selection Evolution, May 2022





On real data from growing pig I

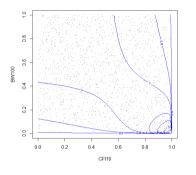




Figure: Contours plot of the fitted copula from pseudo-observations for CFI10 and BW100 (rotated 90 degrees Clayton copula): (left) uniform scale, (middle) gaussian scale, (right) Large-white pigs!





Maximum likelihood estimation

Because the BVs a_i are unobservable, the log-density is

$$\log f_{\mathbf{Y}}(\mathbf{y}) = \log \int_{2N} f_{Y|\mathbf{a}}(\mathbf{y}|\mathbf{a}) f_{\mathbf{a}}(\mathbf{a}) d\mathbf{a}.$$

The MLE is solution with respect to $\xi = (\beta_1, \beta_2, \sigma_1^2, \sigma_2^2, \theta, \sigma_{a_1}^2, \sigma_{a_2}^2, \sigma_{a_{12}})$ of the system equations

$$\frac{\partial}{\partial \xi} \log f_{\mathbf{Y}}(\mathbf{y}; \xi) = 0.$$

Remember that

$$egin{aligned} oldsymbol{a} &= (oldsymbol{a}_1, oldsymbol{a}_2) &\sim \mathcal{N}_{2N}\left(oldsymbol{0}, oldsymbol{G} \otimes A
ight); \ Y_{ij} | Z_i oldsymbol{a}_j &\sim \mathcal{N}\left(Z_i oldsymbol{a}_j + oldsymbol{x}_{ij}eta_j, \sigma_j^2
ight) \ (Y_{i,1}, Y_{i,2}) | (Z_i oldsymbol{a}_1, Z_i oldsymbol{a}_2) & ext{has copula } C_{ heta}; \end{aligned}$$





Stochastic gradient descent algorithm

Remember we have to solve

$$\frac{\partial}{\partial \xi} \log f_{\mathbf{Y}}(\mathbf{y}; \xi) = 0.$$

First note that he Fisher's identity states

$$\frac{\partial}{\partial \xi} \log f_{\mathbf{Y}}(\mathbf{y}; \xi) = \mathbb{E}_{\mathbf{a}|\mathbf{y}} \left(\frac{\partial}{\partial \xi} \log f_{(\mathbf{Y}, \mathbf{a})}(\mathbf{y}, \mathbf{a}; \xi) \right).$$





Stochastic gradient descent algorithm

Remember we have to solve

$$\frac{\partial}{\partial \xi} \log f_{\mathbf{Y}}(\mathbf{y}; \xi) = 0.$$

First note that he Fisher's identity states

$$\frac{\partial}{\partial \xi} \log f_{\mathbf{Y}}(\mathbf{y}; \xi) = \mathbb{E}_{\mathbf{a}|\mathbf{y}} \left(\frac{\partial}{\partial \xi} \log f_{(\mathbf{Y}, \mathbf{a})}(\mathbf{y}, \mathbf{a}; \xi) \right).$$

An SGD algorithm is for $m \in 1, ..., M$ do:

- ▶ Simulate $a^{(m)}$ from the conditional distribution of $a \mid y$
- Update the parameter:

$$\boldsymbol{\xi}^{(m)} = \boldsymbol{\xi}^{(m-1)} + \gamma_m \frac{\partial}{\partial \boldsymbol{\xi}} \log f_{(\boldsymbol{Y},\boldsymbol{a})}(\boldsymbol{y},\boldsymbol{a}^{(m)};\boldsymbol{\xi}^{(m-1)})$$

for a well-chosen learning rate γ_m .



Simplifications

The update step can be rewritten as

$$G^{(m)} = G^{(m-1)} + \gamma_{1,m} \frac{\partial}{\partial \xi} \log f_{\mathbf{a}}(\mathbf{a}^{(m)}),$$

$$(\boldsymbol{\beta}, \sigma_j^2)^{(m)} = (\boldsymbol{\beta}, \sigma_j^2)^{(m-1)} + \gamma_{2,m} \frac{\partial}{\partial \xi} \log f_{\boldsymbol{Y}|\boldsymbol{a}}(\boldsymbol{y}|z^T \boldsymbol{a}^{(m)})$$

and

$$\theta^{(m)} = \theta^{(m-1)} + \gamma_{3,m} \frac{\partial}{\partial \xi} \log f_{Y|a}(y|z^T a^{(m)}).$$





The update of the covariance matrix for the genetic effects is

$$G^{(m)} = G^{(m-1)} + \gamma_{1,m} \frac{\partial}{\partial \xi} \log f_{\mathbf{a}}(\mathbf{a}^{(m)}),$$

A is a very huge and dense matrix, working with A can be numerically complex. But A^{-1} is very sparse! With some simplifications, we can work only with A^{-1} :

$$\begin{split} &\frac{\partial}{\partial \boldsymbol{\xi}} \log f\left(\boldsymbol{a^{(m)}}\right) \\ &= \frac{1}{2} \left(\operatorname{trace}\left(\left(\boldsymbol{G^{(m-1)}} \otimes \boldsymbol{A}\right) \times \left(\nabla \boldsymbol{G^{-1(m-1)}} \right) \otimes \boldsymbol{A^{-1}} \right) - \left(\boldsymbol{a^{(m)}}\right)^T \left(\left(\nabla \boldsymbol{G^{-1(m-1)}} \right) \otimes \boldsymbol{A^{-1}} \right) \boldsymbol{a^{(m)}} \\ &= \frac{1}{2} \left(\boldsymbol{N} \times \operatorname{trace}\left(\boldsymbol{G^{(m-1)}} \times \left(\nabla \boldsymbol{G^{-1(m-1)}} \right) \right) - \left(\boldsymbol{a^{(m)}}\right)^T \left(\left(\nabla \boldsymbol{G^{-1(m-1)}} \right) \otimes \boldsymbol{A^{-1}} \right) \boldsymbol{a^{(m)}} \right) \end{split}$$





It can be rewritten as can be rewritten as

$$G^{(m)} = G^{(m-1)} + \gamma_{1,m} \frac{\partial}{\partial \xi} \log f_{\mathbf{a}}(\mathbf{a}^{(m)}),$$

→ based on trace of sparse matrix

$$(\boldsymbol{\beta}, \sigma_j^2)^{(m)} = (\boldsymbol{\beta}, \sigma_j^2)^{(m-1)} + \gamma_{2,m} \frac{\partial}{\partial \xi} \log f_{\boldsymbol{Y}|\boldsymbol{a}}(\boldsymbol{y}|\boldsymbol{z}^T \boldsymbol{a}^{(m)})$$

and

$$\theta^{(m)} = \theta^{(m-1)} + \gamma_{3,m} \frac{\partial}{\partial \xi} \log f_{Y|a}(y|z^T a^{(m)}).$$

 \rightarrow more complex analytic formulations (derivatives of copula density) but no real challenge to efficiently compute it.



Simulation of the conditional distribution of BVs given observations

- ▶ for Gaussian copula, $a \mid Y$ has explicit Gaussian distribution with covariance $((ZG \otimes AZ^T)^{-1} + (R \otimes I_n)^{-1})^{-1}$.
- → That can be sampled using a Cholesky decomposition



Simulation of the conditional distribution of BVs given observations

- ▶ for Gaussian copula, $a \mid Y$ has explicit Gaussian distribution with covariance $((ZG \otimes AZ^T)^{-1} + (R \otimes I_n)^{-1})^{-1}$.
- → That can be sampled using a Cholesky decomposition
- for non-Gaussian copula, $\mathbf{a}|\mathbf{Y}$ does not have an explicit distribution $(\propto f(Y|\mathbf{a})f(\mathbf{a}))$, and there is not easy resampling sheme.
- → Hybrid MCMC-Metropolips-Gibbs block sampling



Rohmer, T., Bruning, V. and Kuhn, Estelle

G+E copula model to improve the estimation of the genetic parameters in bivariate mixed model, submitted, 2025



On simulations

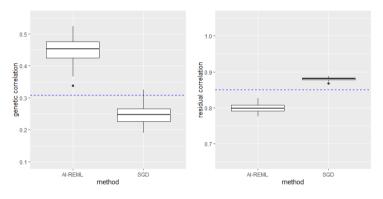
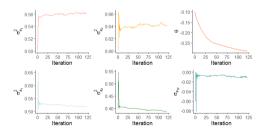


Figure: Boxplot of the estimated genetic correlation (left), residual correlation (right), using B=50 runs. The Clayton model using MLE (SGD) and the Gaussian model using REML are compared. For the two models, data are sampled through Clayton copula for the residual part from 9900 animals under selection using a truncation selection process. Blue dotted line is the true value of the parameter

On real data II

varcomp	$\sigma_{a_1}^2$	$\sigma_{a_2}^2$	$\sigma_{a_{12}}$	$\sigma_{e_1}^2$	$\sigma_{e_2}^2$	heta	h_{CFI10}^2	h_{BW100}^{2}	$ ho_{e}$	iterations
AI-REML	0.56	0.37	-0.02	0.38	$0.\bar{5}4$	-0.18	0.59	0.41	-0.39	7
rC-SGD	0.52	0.39	-0.01	0.56	0.64	-0.28	0.48	0.38	-0.20	121

Table: Estimation of the variance components using Gaussian inference model with Al-REML procedure and using rotated 90 degree Clayton inference model with SGD procedure, using 3 generations, n = 1749, N = 4653.





Next steps

- 1. From bivariate to multivariate analysis.
- $\,\hookrightarrow\,$ Preconditionning the GD step by Fisher's information.



Next steps

- 1. From bivariate to multivariate analysis.
- → Preconditionning the GD step by Fisher's information.
- 2. R package
- \hookrightarrow for users



Next steps

- 1. From bivariate to multivariate analysis.
- → Preconditionning the GD step by Fisher's information.
- 2. R package
- \hookrightarrow for users
- 3. Robustness to a model misspecification? (e.g. copula, or marginal distribution)

