



HAL
open science

Copula Integration for Genetic Parameter Estimation in Bivariate Linear Mixed Models

Victoria Bruning, Estelle Kuhn, Tom Rohmer

► **To cite this version:**

Victoria Bruning, Estelle Kuhn, Tom Rohmer. Copula Integration for Genetic Parameter Estimation in Bivariate Linear Mixed Models. Journées Statistiques du Sud 2024, Jun 2024, Toulouse, France. hal-04633406

HAL Id: hal-04633406

<https://hal.inrae.fr/hal-04633406>

Submitted on 3 Jul 2024

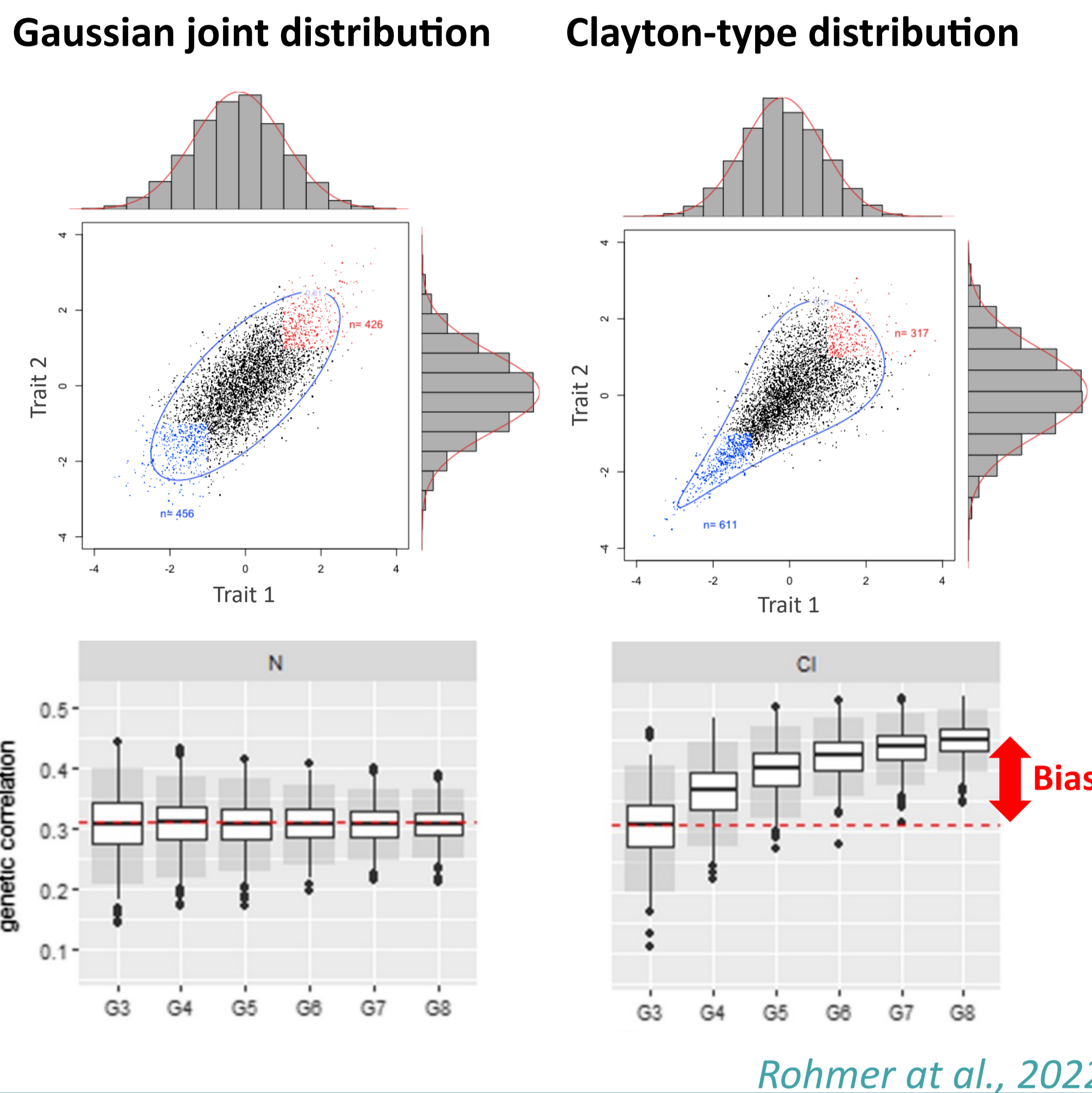
HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

In the bivariate genetic animal model, the estimation methods assume that the traits are gaussian, but in practice, the joint distribution of the two traits can be non-Gaussian. This induces a bias in the genetic parameters estimation when the reproducers are non-randomly selected. The aim of this project was to integrate copula functions, characterising the joint distribution, in the estimation of genetic parameters, and so, reduce this bias.

The problem

Induction of a bias in the REML estimation of the genetic parameters when the residual joint distribution is non-gaussian and the reproducers are non-randomly selected from generation G3.



The genetic animal model

Each trait is determined by an environmental and genetic factor. The breeding value \mathbf{a} is the average additive effects of genes an individual receives from both parents. Thus, it is the primary component that can be selected and is the main focus.

Bivariate Mixed Model

$$\begin{cases} \mathbf{y}_1 = X_1\beta_1 + Z\mathbf{a}_1 + \mathbf{e}_1 \\ \mathbf{y}_2 = X_2\beta_2 + Z\mathbf{a}_2 + \mathbf{e}_2 \end{cases}$$

Labels: Traits, Environment component, Genetic component, Residual term

$\mathbf{a} = (\mathbf{a}_1, \mathbf{a}_2) \sim \mathcal{N}(0, \Sigma)$ with $\Sigma = \mathbf{G} \otimes \mathbf{A}$ (Kinship matrix)

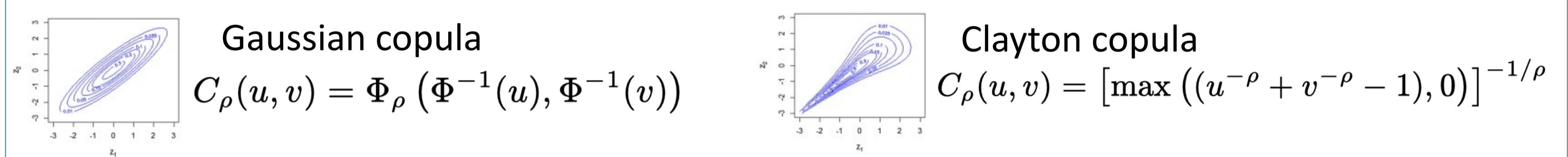
$$\mathbf{G} = \begin{pmatrix} \sigma_{a_1}^2 & \sigma_{a_{12}} \\ \sigma_{a_{12}} & \sigma_{a_2}^2 \end{pmatrix}$$

$$\mathbf{E} = \begin{pmatrix} \sigma_{e_1}^2 & \sigma_{e_{12}} \\ \sigma_{e_{12}} & \sigma_{e_2}^2 \end{pmatrix}$$

Goal: \mathbf{e} is defined by a copula function and Gaussian marginals

The tool : Copula functions

Copulas are functions characterize the joint distribution of the data



The estimation algorithm : Stochastic Gradient with Copulas

The stochastic gradient is an optimization method that updates the parameters using the gradient of the likelihood computed from a simulated conditional distribution

1- Simulate $\mathbf{a}^{(k)} \sim f(\mathbf{a} | \mathbf{y})$, the conditional distribution by Gibbs sampling

2- Maximize $f(\mathbf{y}; \theta)$, the likelihood function

$$\log f(\mathbf{y}; \theta) = \log \int f(\mathbf{y} | \mathbf{a}; \theta) f(\mathbf{a}; \theta) d\mathbf{a}$$



With Gaussian assumption

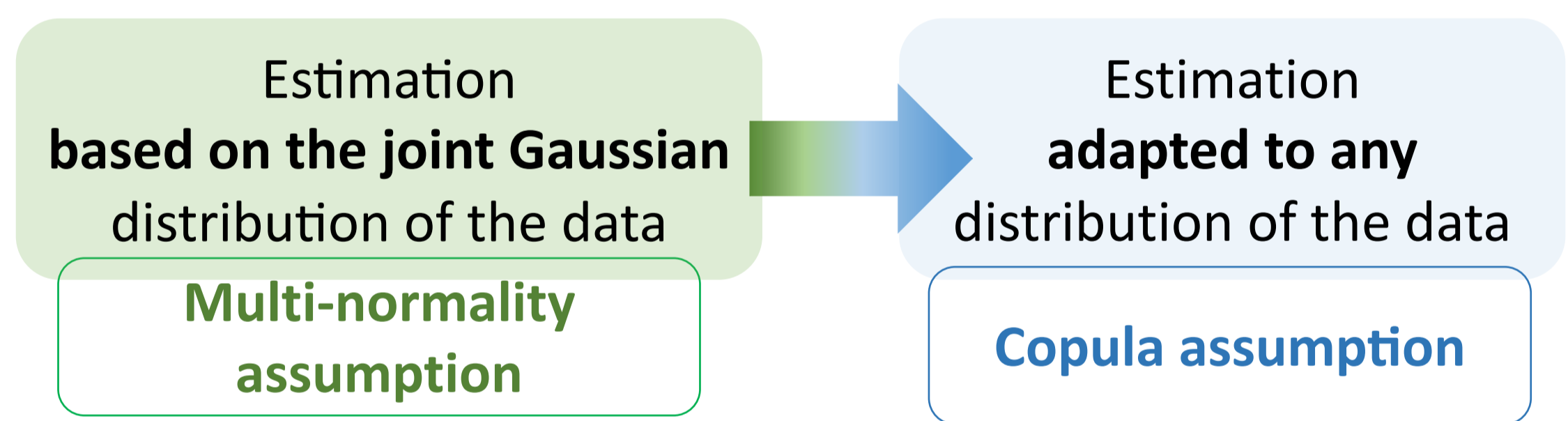
$$f(\mathbf{y} | \mathbf{a}) = \frac{1}{\sqrt{(2\pi)^{2n} \cdot \det(\mathbf{R})}} \exp\left(-\frac{1}{2}(\mathbf{y} - \mathbf{a})^T \mathbf{R}^{-1}(\mathbf{y} - \mathbf{a})\right)$$

Integrating a Copula function

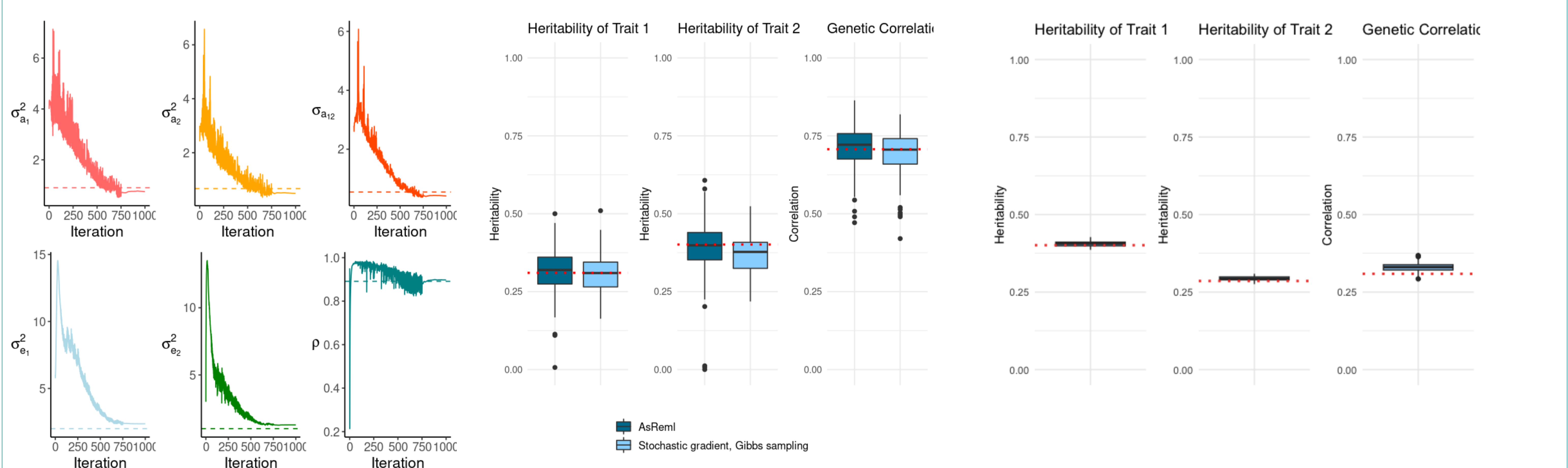
$$f(\mathbf{y} | \mathbf{a}) = \sum_i c(\Phi_1(y_{i1} | \mathbf{a}_{i1}), \Phi_2(y_{i2} | \mathbf{a}_{i2})) \times \prod_{j=1}^2 f_j(y_{ij} | \mathbf{a}_{ij})$$

$$\theta^{(k+1)} = \theta^{(k)} + \text{learning rate} \times \nabla_{\theta} \log f(\mathbf{y}, \mathbf{a}^{(k)}; \theta^{(k)})$$

The goal



Results & Discussion



The simulated data has a Gaussian dependence structure with 720 animals, and the estimation included a Gaussian copula. A- Convergence graph for a single run, B- Boxplot of the estimates obtained for 100 runs by ASReml software and by the developed method. The dotted lines are the true values.

The simulated data has a Clayton-type dependence structure with 7560 animals, and the estimation included a Clayton copula. Boxplot of the estimates obtained for 100 runs by the developed method.

Conclusion

The linear mixed model integrating copulas enables the estimation of genetic parameters when the joint distribution is either Gaussian or defined through another dependence structure, here a Clayton-type copula. The convergence of the stochastic gradient algorithm is observed toward a critical point, and expected toward the MLE which may be closed to the true values for large sample size.

Perspectives

The learning rate will be enriched by a Fisher scoring matrix to standardize the gradient descent direction. This algorithm will be applied on real data. In future work, this could be extended to more than two traits and on data with non-Gaussian marginal distributions.