



# Peer Community Journal

Section: Mathematical & Computational Biology

Research article

Published  
2025-01-06

Cite as

Michel Turbet Delof, Pierre Rivière, Julie C Dawson, Arnaud Gauffreteau, Isabelle Goldringer, Gaëlle van Frank and Olivier David (2025) *Bayesian joint-regression analysis of unbalanced series of on-farm trials*, Peer Community Journal, 5: e4.

Correspondence

[michel.turbet\\_delof@cirad.fr](mailto:michel.turbet_delof@cirad.fr)

Peer-review

Peer reviewed and recommended by PCI Mathematical & Computational Biology, <https://doi.org/10.24072/pci.mcb.100272>



This article is licensed under the Creative Commons Attribution 4.0 License.

## Bayesian joint-regression analysis of unbalanced series of on-farm trials

Michel Turbet Delof<sup>1</sup>, Pierre Rivière<sup>2</sup>, Julie C Dawson<sup>3</sup>, Arnaud Gauffreteau<sup>4</sup>, Isabelle Goldringer<sup>1</sup>, Gaëlle van Frank<sup>1</sup>, and Olivier David<sup>5</sup>

Volume 5 (2025), article e4

<https://doi.org/10.24072/pcjournal.495>

### Abstract

Participatory plant breeding (PPB) is aimed at developing varieties adapted to agroecologically-based systems. In PPB, selection is decentralized in the target environments, and relies on collaboration between farmers, farmers' organisations and researchers. By doing so, evaluation of new genotypes takes genotype x environment (GxE) interactions into account to select for specific adaptation. In many cases, there is little overlap among genotypes assessed from farm to farm because the farmers participating in a PPB project choose which ones to assess on their farm. In addition, on-farm trials can often generate more extreme observations than trials carried out on research stations. These features make the estimation of genotype, environment and interaction effects more difficult. This challenge is not unique to PPB, as many breeding programs use sparse testing or incomplete block designs to evaluate more genotypes, however in PPB genotypes are not always assigned randomly to environments. To explore methods of overcoming these challenges, this article tests various data analysis scenarios using a Bayesian approach with different models and a real wheat PPB dataset over 11 years. Four morpho-agronomic traits were studied, representing over 1000 GxE combinations from 189 on-farm trials. This dataset was severely unbalanced with more than 90% of GxE combinations missing. We compared various Bayesian Finlay-Wilkinson models and found that placing hierarchical distributions on model parameters and modelling residuals using a Student's t distribution jointly improved the estimates of main effects and interactions. Environment effects were the most important and explained more than 50% of the variance of observations. This statistical framework allowed us to estimate two indicators of genotype stability (one static and one dynamic) despite the high disequilibrium of the data. We found differences in mean and stability between genotype categories, with registered varieties consistently shorter (-30 cm) and containing less protein (-0.3%) than other types of varieties. The methods developed could be used for evaluation and/or selection within networks of various stakeholders such as farmers, gardeners, plant breeders or managers of genetic resource centres.

<sup>1</sup>UMR GQE-Le Moulon, Université Paris-Saclay - INRAE - CNRS - AgroParisTech, 91190, Gif-sur-Yvette, France, <sup>2</sup>Métis, 47213, Prayssas, France, <sup>3</sup>Department of Plant and Agroecosystem Sciences, University of Wisconsin-Madison, WI 53706, Madison, USA, <sup>4</sup>UMR Agronomie, Université Paris-Saclay - AgroParisTech - INRAE, 91120, Palaiseau, France, <sup>5</sup>Université Paris-Saclay, INRAE, MaIAGE, 78350, Jouy-en-Josas, France

Peer Community Journal is a member of the  
Centre Mersenne for Open Scientific Publishing  
<http://www.centre-mersenne.org/>

e-ISSN 2804-3871

## 1. Introduction

Developing new varieties adapted to Organic Agriculture (OA), agroecological and low input systems is a major concern to achieve improvements in agricultural sustainability (Wolfe et al., 2008). In OA, the use of synthetic inputs (nitrogen, phytochemicals) is not allowed, therefore, cropping environments are not standardized by inputs and varieties grow in more diverse conditions from farm to farm (Dawson et al., 2008). These environments are more sensitive to pedoclimatic conditions, yearly weather, farmers' management practices and interactions between these factors (Desclaux et al., 2008).

In order to develop varieties adapted to such a diversity of environments two strategies can be used: (i) centralized and indirect selection, or (ii) decentralized and direct selection. The key difference between these approaches is the way they take genotype-by-environment ( $G \times E$ ) interactions into account. These interactions are considered by plant breeders as the main factor limiting the efficiency of the response to selection in breeding programs (Ceccarelli et al., 2001). In centralized and indirect selection, breeding lines are evaluated and selected at a few research stations assumed to represent the target environments. This is efficient if there is a high additive genetic correlation between the trait measured on the station and the same trait measured in the target environment, and if the narrow sense heritability is high in the selection environment (Falconer, 1960).

Decentralized selection can take account of  $G \times E$  interactions that are important in OA (Dawson et al., 2008; Murphy et al., 2007). In this approach, the selection and evaluation environments are very close to the target environments (the production environments of farms). Selection then maximizes the use of the reproducible part of  $G \times E$  interactions to select for specific adaptations (Annicchiarico et al., 2010). This method is close to direct selection and has been shown to be effective (Annicchiarico et al., 2010; Ceccarelli et al., 2001; Murphy et al., 2007; Smith et al., 2001; Virk et al., 2005).

Many participatory plant breeding (PPB) programs have been carried out over the last 20 years targeting low-input farming systems in the Global South and also OA and agroecological systems in Europe and North America (Ceccarelli and Grando, 2020). A few programs tested different experimental designs and specific statistical methods to analyze data taking  $G \times E$  into account (Mohammadi et al., 2011; Snapp and Silim, 2002). Recently, participatory variety trials using crowdsourcing have been used in several countries with great success (Van Etten et al., 2019). These methods typically use an experimental design called a triadic comparison of technologies (tricot), followed by an analysis of variety ranks (Beza et al., 2017). In the tricot design, large numbers of farmers each compare three variety subsets from the complete set of entries, and provide direct comparison rankings among them for a few traits (i.e. best/middle/worst). By using ranking methods and structuring the entry distribution as an incomplete block design, this allows for comparisons of larger numbers of varieties without overburdening individual farmers. These design options enhance breeders' ability to engage farmers in trialing experimental lines, since on-farm trials are often limited by space and farmers' time. Trialing a few experimental lines, including a check line or variety that is replicated across sites is more realistic for farmers than implementing a fully replicated design. Triadic methods are very useful in many situations, but they are not applicable to more mature farmer-breeder networks, where the choice of varieties and cropping practices is made by farmers. In addition, farmers may want to test different numbers of varieties, with some testing just a few and others several dozen. Farmers also wish to have access to quantitative data rather than simple rankings, so a non-parametric ranking of varieties without assumptions about distribution will not produce a satisfactory analysis for this purpose.

One program with such concerns is a wheat PPB program that started in France in 2005, as a collaboration between INRAE GQE-Le Moulon and the Farmers' Seed Network (Réseau Semences Paysannes, RSP). This PPB program had three objectives: (i) develop varieties adapted to farmers' practices and needs (organic management, artisanal bread quality ...) using a participatory approach, (ii) develop strategies for preserving genetic diversity through on-farm dynamic

management and breeding, and (iii) learn from and improve farmers' individual and collective breeding methods and diffuse successful methods broadly.

In this program, farmers conducted trials with different varieties developed through their own breeding efforts to determine which variety was best suited to their production systems (Turbet Delof, 2024). The research team provided methods to assist farmers in interpreting these trials, aiming to empower them (Rivière et al., 2015a; Turbet Delof, 2024; van Frank, 2018) and to provide general knowledge about these varieties (Goldringer et al., 2020; Rivière et al., 2015b; van Frank et al., 2020). When farmers seek to incorporate and evaluate new populations in their trials, they often struggle with a lack of information on which populations to select. This highlights the need for support in varietal choice, including information on the average performance and stability of varieties within the trial network. Specifically, interannual stability is crucial as it relates to both agronomic and economic risks. Static stability describes the response of a genotype that maintains a constant performance across environments, while dynamic stability describes the response of a genotype showing a constant difference with an environmental reference (generally the average response of all the genotypes, Annicchiarico, 2002).

As very few varieties were common to all the trials and many varieties were tested in a limited number of trials, the resulting series of trials was very unbalanced, so that the estimation of variety average performances and stabilities was difficult. Joint regression is a robust method for estimating genetic main effects and stability with incomplete datasets (Finlay and Wilkinson, 1963; Pereira et al., 2007; Yates and Cochran, 1938). It is based on the Finlay-Wilkinson (FW) model, which is parsimonious since the interaction effect between a genotype and an environment is modelled as the product of a genotype stability parameter, called sensitivity, and the environment main effect. Various Finlay-Wilkinson models have been used in a frequentist framework, in which environment effects were either fixed or random (Nabugoomu et al., 1999; Ng and Williams, 2001; Patterson and Silvey, 1980). In the latter case, environment effects were assumed to come from a common distribution, thereby leading to shrunk estimates. FW models in which genetic main effects, environment main effects and genetic sensitivities (FW coefficient of regression) are all random effects have recently been developed. These have been implemented in a Bayesian framework and when they include random effects, these are called hierarchical models (Carlin and Louis, 2009; Robert, 2007). Thus far, these models have been used to analyze slightly unbalanced trials (Lian and Campos, 2016). Hierarchical joint regression has also been used to analyze very unbalanced simulated data (van Frank et al., 2019). This simulation study has shown that genotypes should be tested in sufficiently many trials in order to estimate their main effects and sensitivities reliably. However, this method had not been used to analyze real and very unbalanced trials. Thus, it was not clear if it could cope with the actual levels of unbalanced data seen in the French PPB on-farm trials and what insight it could give into the behavior of genotypes across environments.

Extreme data is an important issue in data analysis. In multi-environment trials (MET), they may come from either (1) errors between scoring and data formatting (measurement error, wrong labelling, etc.), or (2) environmental heterogeneity in the trial (weed infestation, soil fertility, etc.), or (3) the heterogeneity of the responses of the varieties tested between trials ( $G \times E$  interaction). In our PPB program, as cultivation environments are less controlled, extreme observations (types 2 and 3) could be more frequent than expected. This could reduce the precision of estimates based on the normal distribution. Extreme observations could be removed from the dataset to solve this problem, but it is difficult to decide which observations to remove. If too many extreme observations are removed, then the variability of the data may be underestimated and the precision of the statistical analysis overestimated. Alternatively, statistical methods that are robust to extreme observations may be used (Hampel et al., 2011; Huber and Ronchetti, 1981). Various robust methods have been developed in a frequentist or a Bayesian framework, in particular methods consisting in replacing the normal distribution by a Student's  $t$  distribution in statistical models. This distribution is more robust to extreme observations than the normal distribution, because it has heavier tails (Carlin and Polson, 1991; Choy and Chan, 2008; Lange et al., 1989; Rosa et al., 2003). It has been used to handle the extreme observations of a single

trial in a Bayesian framework (Besag and Higdon, 1999; Cao et al., 2022; Gianola et al., 2018). However, to our knowledge, it has not been used to analyze an unbalanced network of trials.

This study was aimed at developing statistical methods for analyzing series of on-farm trials, and at improving the assessment of varieties of the wheat PPB program by using the information at the level of the network. As our dataset was very unbalanced and could include extreme observations, we compared several Finlay-Wilkinson models, in particular hierarchical models and models based on the  $t$  distribution. These models were developed in a Bayesian framework, since this framework is rigorous and since it facilitates the implementation of complex models (Carlin and Louis, 2009; Robert, 2007). Finally, the best Finlay-Wilkinson model we obtained was used to analyze our data and characterize the behaviour of our varieties across environments.

## 2. Materials and methods

In our study, a population variety is defined as a set of individuals which may be different but which are derived from certain agronomic practices, and a germplasm as any biological entity whose individuals are derived from the same breeding process, including varieties registered in the official catalog, landraces, historic varieties, mixtures or populations stemming from crosses. An environment is the combination of a farm and a year. The main notations used in this study are shown in Tab.1.

**Table 1** – Main notations.

Notation	Meaning
PPB	Participatory plant breeding
OA	Organic agriculture
RSP	Réseau Semences Paysannes, French farmers' seed network
MET	Multi-environment trial
$G \times E$	Genotype $\times$ environment interaction
FW	Finlay Wilkinson
MCMC	Markov chain Monte Carlo
$\alpha$	Germplasm main effect
$\theta$	Environment main effect
$\eta$	Germplasm sensitivity (FW coefficient)
$S^2$	Germplasm static stability
$W$	Germplasm ecovalence (a dynamic stability)
LOO	Leave one out
$\text{elpd}_{\text{loo}}$	LOO expected logarithmic predictive density

### 2.1. Statistical methods

**2.1.1. Models.** We consider methods for analyzing series of on-farm trials in two steps (Patterson, 1997; Patterson and Silvey, 1980). First, germplasm means are estimated using within-trial analyses, taking into account any block effects (spatial effects). Then, these estimates are analyzed using a between-trial analysis. In the between-trial analysis, the phenotypic value  $Y_{ij} \in \mathbb{R}$  for a given trait  $Y$ , germplasm  $i$  and environment  $j$  is assumed to be equal to

$$Y_{ij} = \mu_{ij} + \varepsilon_{ij},$$

where  $(i, j) \in \mathcal{C}$ ,  $\mathcal{C}$  is the set of the germplasm  $\times$  environment combinations occurring in the dataset,  $\mu_{ij} \in \mathbb{R}$  is an expectation term, and  $\varepsilon_{ij} \in \mathbb{R}$  is a between-trial residual term.

In models ADHs and ADHn, the expectation term is modelled as additive effects of both the germplasm and the environment without interaction:

$$\mu_{ij} = \alpha_i + \theta_j,$$

where  $\alpha_i \in \mathbb{R}$  is the main effect of germplasm  $i$ , and  $\theta_j \in \mathbb{R}$  is the main effect of environment  $j$ . Models FWHs, FWs and FWHn model  $G \times E$  interactions using the Finlay-Wilkinson regression,

**Table 2** – The five models compared.

Model	Expectation term	Residual term	Prior distribution
ADHn	Additive	Normal	Hierarchical
ADHs	Additive	Student	Hierarchical
FWHn	Finlay Wilkinson	Normal	Hierarchical
FWHs	Finlay Wilkinson	Student	Hierarchical
FWs	Finlay Wilkinson	Student	Weakly informative

also called joint-regression, model (Finlay and Wilkinson, 1963; Yates and Cochran, 1938). In these models, the expectation term is assumed to be equal to

$$\mu_{ij} = \alpha_i + \theta_j + \eta_i\theta_j,$$

where  $\eta_i \in \mathbb{R}$  is the sensitivity of germplasm  $i$  to environments (regression coefficient, Perkins and Jinks, 1968). As the average sensitivity is equal to 0, a germplasm with  $\eta_i > 0$  is more sensitive and a germplasm with  $\eta_i < 0$  is less sensitive to environments than a germplasm with the average sensitivity. In these models, a part of the interaction between germplasm  $i$  and environment  $j$  is modelled as a multiplicative term  $\eta_i\theta_j$ . The Finlay-Wilkinson coefficient is considered as both a static and a dynamic indicator of stability (Becker and Leon, 1988; Lin et al., 1986). In this model, statically stable genotypes have a coefficient close to -1. Dynamically stable genotypes have a coefficient close to zero, but having a coefficient close to zero is not sufficient to determine dynamic stability, this also depends on the amount of  $G \times E$  variation that remains unexplained by the model.

We consider series on-farm trials where most of the germplasm are not replicated within the trials. For such trials, the standard errors of germplasm means provided by the within-trial analyses are not precise. Thus, these standard errors are not taken into account, and the between-trial residuals are assumed to be homoscedastic (Patterson, 1997; Patterson and Silvey, 1980). In models ADHn and FWHn, the distribution of these residuals is assumed to be normal:

$$\varepsilon_{ij} \sim \mathcal{N}(0, \sigma_\varepsilon^2),$$

where  $\mathcal{N}(0, \sigma_\varepsilon^2)$  is the normal distribution with expectation 0 and variance  $\sigma_\varepsilon^2$ . However, to limit the influence of extreme values on the results of the analyses, we also consider models based on Student's  $t$  distributions. Thus, in models FWHs, FWs and ADHs, the distribution of the error term is assumed to be equal to

$$\varepsilon_{ij} \sim t(0, \sigma_\varepsilon^2, \nu),$$

where  $t(0, \sigma_\varepsilon^2, \nu)$  is the Student's  $t$  distribution with dispersion parameter  $\sigma_\varepsilon^2 > 0$  and  $\nu > 2$  degrees of freedom. We assume that  $\nu > 2$  to ensure that the expectation and the variance of  $\varepsilon_{ij}$  are defined and finite. In models FWHs, FWs and ADHs, the variance of  $\varepsilon_{ij}$  is equal to  $\nu\sigma_\varepsilon^2/(\nu-2)$ . The normal distribution can be considered as a  $t$  distribution with  $\nu$  tending to  $+\infty$ . For additive models, the between-trial residuals combine the  $G \times E$  effects and within-trial errors, i.e. experimental errors and environmental heterogeneity in each trial, while for FW models, they combine the part of  $G \times E$  effects not explained by the multiplicative term  $\eta_i\theta_j$  and within-trial errors. Student residuals better handle data heterogeneity than normal residuals, since they can be written as (Simar, 2002)

$$\varepsilon_{ij} \sim \mathcal{N}(0, \sigma_{ij}^2), \quad \sigma_{ij}^{-2} \sim \Gamma(\nu/2, \nu\sigma_\varepsilon^2/2),$$

where  $\Gamma(\nu/2, \nu\sigma_\varepsilon^2/2)$  is the gamma distribution with shape parameter  $\nu/2$  and rate parameter  $\nu\sigma_\varepsilon^2/2$ .

**2.1.2. Prior distribution.** The statistical methods are implemented in a Bayesian framework, so that a joint prior distribution is placed on model parameters. Weakly informative prior distributions are placed on  $\sigma_\varepsilon$  and  $\nu$  (Cao et al., 2022; Gelman, 2006; Juárez and Steel, 2010):

$$\sigma_\varepsilon \sim \mathcal{N}^+(0, \lambda_\varepsilon^2), \quad \nu = 2 + \gamma, \quad \gamma \sim \Gamma(2, 0.1),$$

where  $\lambda_\varepsilon$  is a known prior value of the standard deviation of the trait, and  $\mathcal{N}^+(0, \lambda_\varepsilon^2)$  is the normal distribution restricted to positive values with parameters 0 and  $\lambda_\varepsilon^2$ .

Since series of on-farm trials are often unbalanced and often involve many germplasm and environments,  $\alpha_i$ ,  $\theta_j$  and when present  $\eta_i$  are assumed to follow hierarchical distributions in all the models except model FWs:

$$\alpha_i \sim \mathcal{N}(\mu_Y, \sigma_\alpha^2), \quad \eta_i \sim \mathcal{N}(0, \sigma_\eta^2), \quad \theta_j \sim \mathcal{N}(0, \sigma_\theta^2),$$

where  $\mu_Y$ ,  $\sigma_\alpha$ ,  $\sigma_\eta$  and  $\sigma_\theta$  are unknown parameters. Then, weakly informative prior distributions are placed on the hyperparameters  $\mu_Y$ ,  $\sigma_\alpha$ ,  $\sigma_\eta$  and  $\sigma_\theta$ :

$$\mu_Y \sim \mathcal{N}(\lambda_\mu, \lambda_\varepsilon^2), \quad \sigma_\alpha \sim \mathcal{N}^+(0, \lambda_\varepsilon^2), \quad \sigma_\theta \sim \mathcal{N}^+(0, \lambda_\varepsilon^2), \quad \sigma_\eta \sim \mathcal{N}^+(0, 0.75^2),$$

where  $\lambda_\mu$  is a known prior value of the trait mean. Germplasm main effects, environment main effects, germplasm sensitivities and residuals are assumed to be independent given the hyperparameters,  $\sigma_\varepsilon$  and  $\nu$ . In model FWs, the hierarchical distributions of  $\alpha_i$ ,  $\eta_i$  and  $\theta_j$  are replaced by weakly informative prior distributions:

$$\alpha_i \sim \mathcal{N}(\mu_Y, \lambda_\varepsilon^2), \quad \eta_i \sim \mathcal{N}(0, 0.75^2), \quad \theta_j \sim \mathcal{N}(0, \lambda_\varepsilon^2).$$

The values chosen for  $\lambda_\varepsilon$  and  $\lambda_\mu$  are in Appendix A.1.

In conclusion, five models are considered, which model the expectation term, the residual term and the prior distribution differently (Tab. 2). The main model of interest is FWs, the other models being mainly used for assessing model FWs.

**2.1.3. Posterior distribution.** Bayesian inference is based on the posterior distribution of model parameters. This distribution is estimated using Markov chain and Monte Carlo (MCMC) methods. These methods simulate the values of model parameters according to a Markov chain that converges to the posterior distribution of these parameters (Robert, 2007). They are implemented using R (R Core Team, 2024) and the package `rstan` (Stan Development Team, 2024), that performs Hamiltonian Monte Carlo (HMC) sampling. This method aims at reducing the correlation between successive sampled values by using a proposal distribution based on Hamiltonian dynamics (Neal, 2011).

**2.1.4. Model comparison.** The predictive ability of models is compared using leave-one-out cross-validation, which seems more appropriate than Bayes factors for selecting models that approximate the process generating the data (Lartillot, 2023). We estimate the expected logarithmic predictive density using the R package `LDD` (Vehtari et al., 2017). This criterion is equal to

$$\text{elpd}_{\text{loo}} = \sum_{(i,j) \in \mathcal{C}} \ln(p(Y_{ij} | Y_{-ij})),$$

where  $Y_{-ij}$  is the dataset without observation  $Y_{ij}$ , and  $p(Y_{ij} | Y_{-ij})$  is the leave-one-out posterior density of  $Y_{ij}$ . The larger this criterion, the better the agreement between the model and the data. This criterion is also used to identify extreme observations. The quantity  $\ln(p(Y_{ij} | Y_{-ij}))$  can be understood as the contribution of observation  $Y_{ij}$  to  $\text{elpd}_{\text{loo}}$ . Observations with low contributions are unlikely and can be considered extreme observations.

For main effects and sensitivities, we estimate the average standard deviation of estimates, which allows us to estimate the precision of the analysis. To be able to compare the precision between traits, for  $\alpha$  and  $\theta$  we estimate the average coefficient of variation by dividing this standard deviation by the general average  $\mu_Y$ .

**2.1.5. Variance decomposition.** In order to assess the importance of model terms, the variance of an observation is decomposed for the main model FWs. Since  $\alpha_i$ ,  $\theta_j$ ,  $\eta_i$  and  $\varepsilon_{ij}$  are conditionally independent, the terms  $\theta_j$  and  $\eta_i\theta_j$  are not correlated, and the variance of an observation given the hyperparameters,  $\sigma_\varepsilon^2$  and  $\nu$  is equal to

$$\text{Var}(Y_{ij}) = \text{Var}(\alpha_i + \theta_j + \eta_i\theta_j + \varepsilon_{ij}) = \sigma_\alpha^2 + \sigma_\theta^2 + \sigma_\eta^2\sigma_\theta^2 + \text{Var}(\varepsilon_{ij}).$$

The variance of  $\varepsilon_{ij}$  is equal to  $\nu\sigma_\varepsilon^2/(\nu-2)$  for model FWHs. The proportions of variance explained by the germplasm main effect, the environment main effect and the interaction effect are equal to

$$\pi(\alpha) = \frac{\sigma_\alpha^2}{\text{Var}(Y_{ij})}, \quad \pi(\theta) = \frac{\sigma_\theta^2}{\text{Var}(Y_{ij})}, \quad \pi(\eta\theta) = \frac{\sigma_\eta^2\sigma_\theta^2}{\text{Var}(Y_{ij})}.$$

$\pi(\alpha)$  is also called broad-sense heritability. The proportion of variance explained by the model (coefficient of determination) is equal to

$$R^2 = \pi(\alpha) + \pi(\theta) + \pi(\eta\theta) = \frac{\sigma_\alpha^2 + \sigma_\theta^2 + \sigma_\eta^2\sigma_\theta^2}{\text{Var}(Y_{ij})}.$$

We also estimate the proportion of the variance of  $G \times E$  interactions and experimental errors that is explained by the multiplicative term  $\eta_i\theta_j$ , defined by

$$\rho = \frac{\text{Var}(\eta_i\theta_j)}{\text{Var}(\eta_i\theta_j + \varepsilon_{ij})} = \frac{\sigma_\eta^2\sigma_\theta^2}{\sigma_\eta^2\sigma_\theta^2 + \text{Var}(\varepsilon_{ij})}.$$

**2.1.6. Characterization of germplasm.** The main effect and sensitivity of each germplasm are estimated using model FWHs. In addition, two stability indicators are estimated for each germplasm, the static stability  $S_i^2$  (Becker and Leon, 1988) and the ecovalence  $W_i$  (Wricke, 1962) which is an indicator of dynamic stability. Due to data imbalance, the empirical estimates of these indicators are biased. Thus, we define stability indicators by means of theoretical variances using model FWHs (Cotes et al., 2006; Piepho, 1999). Using the independence assumptions of the model, we obtain for germplasm  $i$ ,

$$\begin{aligned} W_i &= \text{Var}(\eta_i\theta_j + \varepsilon_{ij}) = \eta_i^2\sigma_\theta^2 + \text{Var}(\varepsilon_{ij}), \\ S_i^2 &= \text{Var}(\theta_j + \eta_i\theta_j + \varepsilon_{ij}) = (1 + \eta_i)^2\sigma_\theta^2 + \text{Var}(\varepsilon_{ij}) = (1 + 2\eta_i)\sigma_\theta^2 + W_i. \end{aligned}$$

The larger these indicators, the less stable the germplasm. Becker (1981) applied the same decomposition with the empirical variances.

We also perform pairwise comparisons between germplasm types (e.g., cross, landrace, registered variety, mixture of germplasm and historic variety). For example, for main effects, we compute the average main effect of type  $k$ , denoted by  $\bar{\alpha}_k$ . The comparison between types  $k$  and  $l$  is considered as significant if the 95% credible interval of  $\bar{\alpha}_k - \bar{\alpha}_l$  does not contain 0. Then, germplasm types are grouped into significantly different sets using these pairwise comparisons and an "insert-and-absorb" algorithm (Piepho, 2004).

## 2.2. Wheat PPB program

**2.2.1. Germplasm.** We studied 206 unique germplasm covering different "germplasm types": 98 "cross" germplasm resulting from crosses made either on the farm or at the research station (Rivière et al., 2015b), 50 "landraces", i.e. population varieties grown before 1884 (date of creation of Dattel, the first wheat variety from a controlled cross), 30 "historic varieties", developed by professional breeding before 1950, 17 "mixtures", which were generally complex, with numerous genotypes from potentially all the other germplasm types. In addition, 11 "registered varieties" after 1950 and widely used in organic farming were included: Maitre Pierre (1954), Poncheau (1956), Renan (1990), Ataro (2004), Pollux (2004), Rubisco (2012), Hendrix (2012), Kampmann selected from Renan, and Hermes (1982), Alauda (2004) and Goldritter (2013), all three selected from Probus (1957).

**2.2.2. Experimental designs.** The data analyzed were collected between 2008 and 2019. The wheat PPB program followed numerous experimental designs due to the different constraints of farmers, collectives and researchers. The designs have been grouped into three classes (Tab. 3). Some experimental designs (without blocks with replicated control, and incomplete blocks with two blocks) were co-designed to be adapted to breeders' objectives, farmers' constraints and agricultural routines (Dawson et al., 2011). In these designs, the germplasm common to all farms (controls) were collectively chosen by farmers and researchers, while each farmer individually chose the additional germplasm to be cultivated in his farm. At the beginning the control

was a selection in a landrace, and after 2014 it was a new variety from a cross selected by farmers. Most of the germplasm were not replicated within the trials. All varieties were randomized within farms, but not randomized between farms. Some designs (complete blocks, incomplete blocks) were used to address specific research questions such as the study of the evolution of traits (Rivière et al., 2015b), local adaptation (van Frank et al., 2020) or the evaluation of agronomic performance (Goldringer et al., 2020). Some unreplicated trials corresponded to trials with replications but for which measurements could not be performed in some replications.

**Table 3** – Experimental designs of the 189 trials used in the statistical analysis. Nb: number, Environment: combination of a year and a farm.

Designs	Nb of blocks	Nb of repeated germplasm	Nb of gemplasm by environment	Nb of environments
Complete blocks	2 to 3	6 to 45	7 to 45	24
Incomplete blocks	2 to 4	3 to 49	6 to 81	31
Without blocks		1 to 22 0	5 to 79	102 32

**2.2.3. Data collected.** Four traits were studied, plant height (60% of the data was the average height of 25 individuals and 40% was the overall height of the microplot, mm), spike weight (mean of 25 individual measures, g), protein content of the grain (on the microplot, measured with NIRS technology at INRAE Clermont-Ferrand France, %) and thousand kernel weight (TKW, measured on the microplot, g). These four traits were among those collectively chosen by farmers and researchers to be measured during the PPB program (Tab. 4). Plant height was measured in the field, while the other traits were measured after harvest at the research station on samples of spikes sent by farmers. Outliers with respect to agronomic knowledge of the traits were excluded (for example, a plant taller than three meters).

van Frank et al. (2019) analyzed the sensitivity of the hierarchical FW model to different MET set-ups with simulated data. They found that, in contrast to the environmental effects, the germplasm effects and FW coefficients were difficult to estimate. This is why they recommended that a large number of environments be used and that the germplasm be repeated sufficiently. We have therefore made a selection of the data and kept the environments with at least five germplasm and the germplasm that were present in at least four environments. Thus, the data analyzed comprised 70 to 76% of the initial data, depending on the trait.

The multi-environment data were very unbalanced, with most of the germplasm occurring in a limited number of environments (the median number of replicates across environments was seven, and about 20% of the germplasm were replicated in four environments only). For each trait, the number of observations was between 1300 and 2000 and the measures were spread over more than nine years (Tab. 4).

These data were analyzed using the models of Tab. 2. As the dataset was very unbalanced, it was not clear if model parameters were identifiable. Thus, for each variable, the identifiability of germplasm main effects and environment main effects was studied for the additive model. We checked that the rank of the design matrix of the model was equal to  $1 + (I - 1) + (J - 1)$ , where  $I$  was the number of germplasm and  $J$  the number of environments (p. 50, Silvey, 1975). For the FW model, identifiability was more difficult to study because the model was nonlinear. Thus, we restricted ourselves to studying local identifiability near an estimate of model parameters (Chap. 2, Walter and Pronzato, 1997). First, a linear approximation of the model was carried out using a Taylor expansion. Then, we checked that the rank of the design matrix of this linear model was equal to  $1 + (I - 1) + (J - 1) + (I - 1)$ .

Four MCMC chains were run independently to test for convergence. The initial values of each chain were taken randomly. For each chain, the burn-in consisted of 200 iterations, then 5,000 iterations were performed for all models, except FWs where 10,000 iterations were required. The average calculation time (for a given trait and a given model) was 6 minutes and the maximum time was 22 minutes, with a computer *intel CORE i7*©. Estimates of the Gelman-Rubin statistic



were smaller than 1.02 and the effective sample size was greater than 400 for each parameter in all tested models.

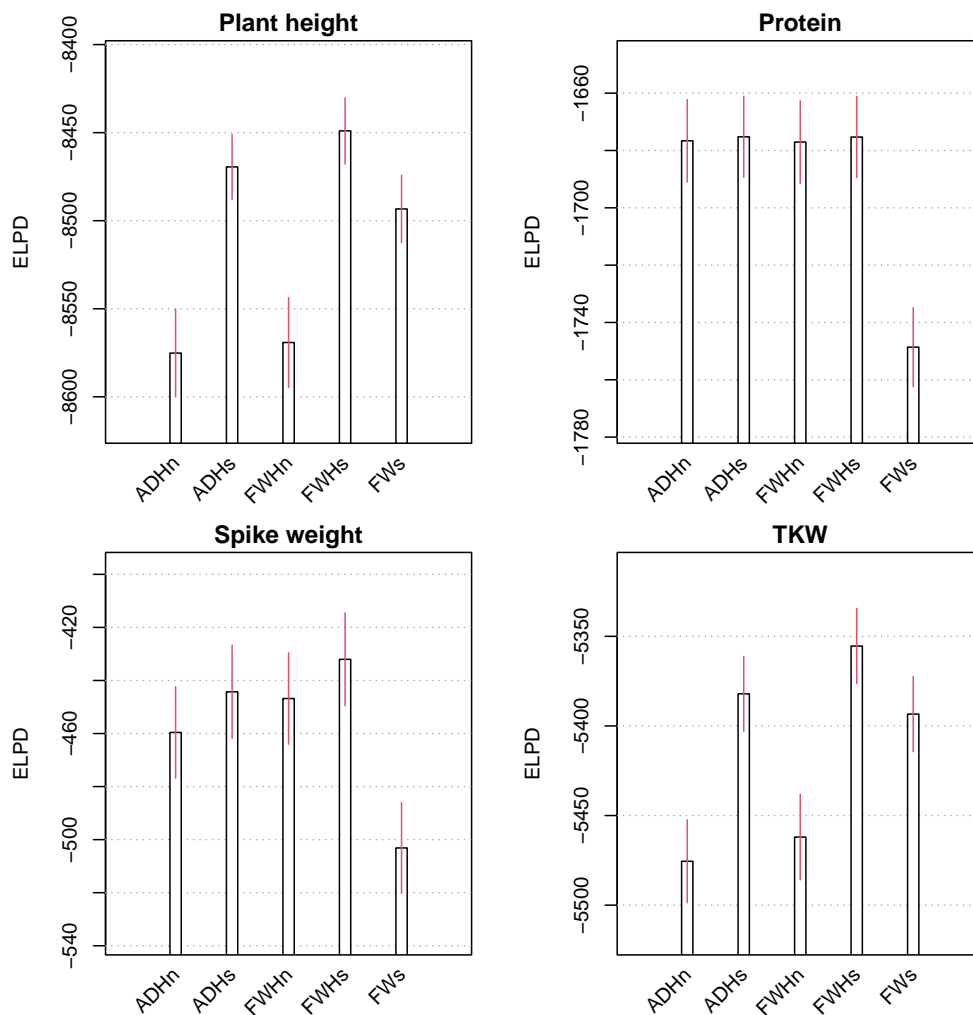
**Table 4** – Description of the dataset. Disequilibrium: percentage of missing values in the Germplasm x Environment table.

Trait	Observations	Germplasm	Environments	Disequilibrium	Farms	Years
Plant height	1437	124	117	90	44	11
Spike weight	1804	172	148	93	52	10
Protein	1332	144	111	92	44	9
TKW	1982	177	165	93	58	11

### 3. Results

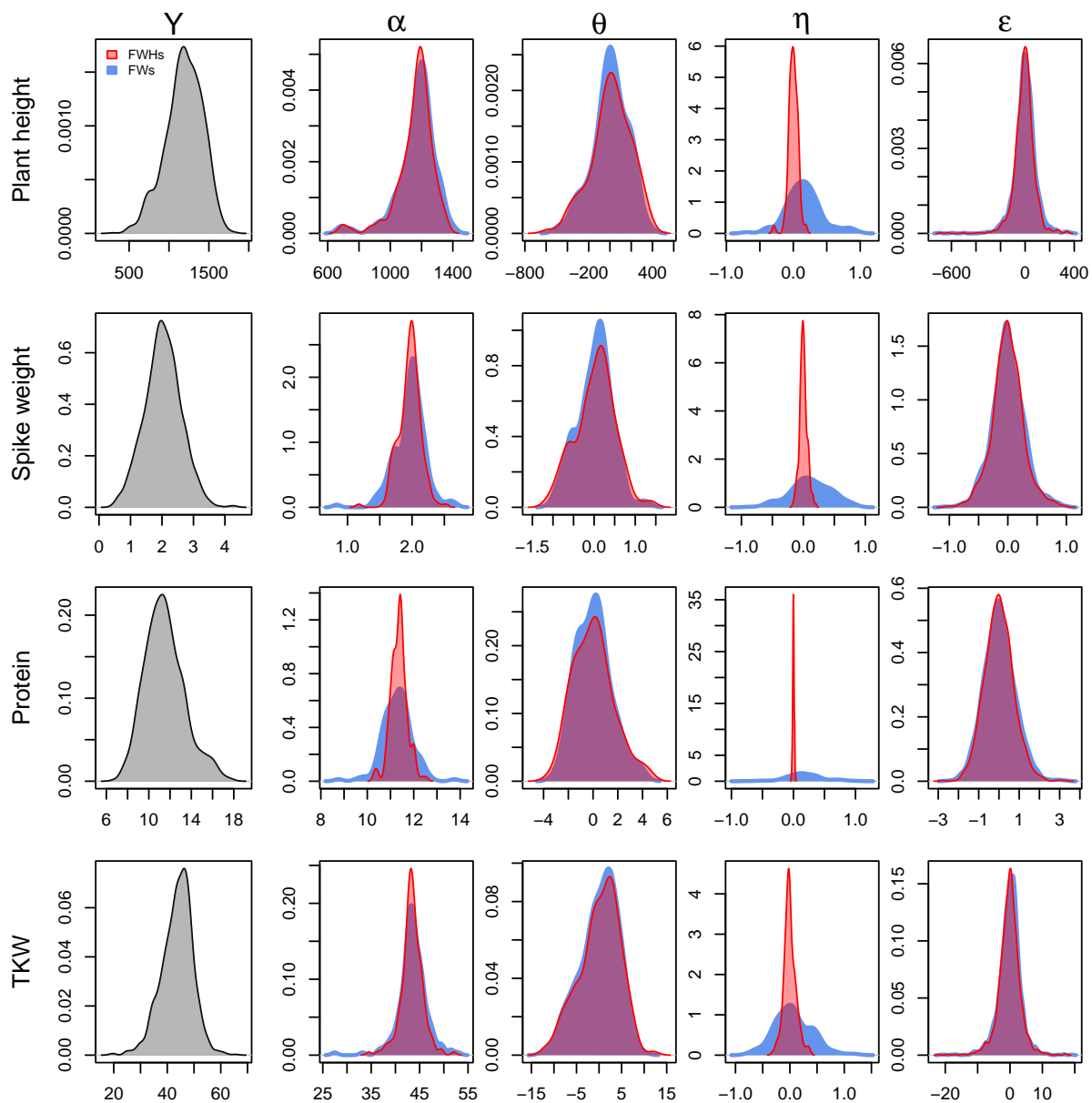
#### 3.1. Predictive capacity of models

According to the  $elpd_{loo}$  criterion, the non-hierarchical FWs model was less predictive than the hierarchical FWHs model for all the traits (Fig. 1). Using the latter model shrank the estimates of  $\eta$  and sometimes  $\alpha$  (Fig. 2).



**Figure 1** – Predictive capacity of models.  $elpd_{loo}$  and its associated standard error for the four studied traits.

With the non-hierarchical model (FWs), some estimates ( $\alpha_i$  and  $\eta_i$ ) seemed to be unreliable, in particular some germplasm means were extreme and some FW coefficients were larger than 1 or smaller than -1.

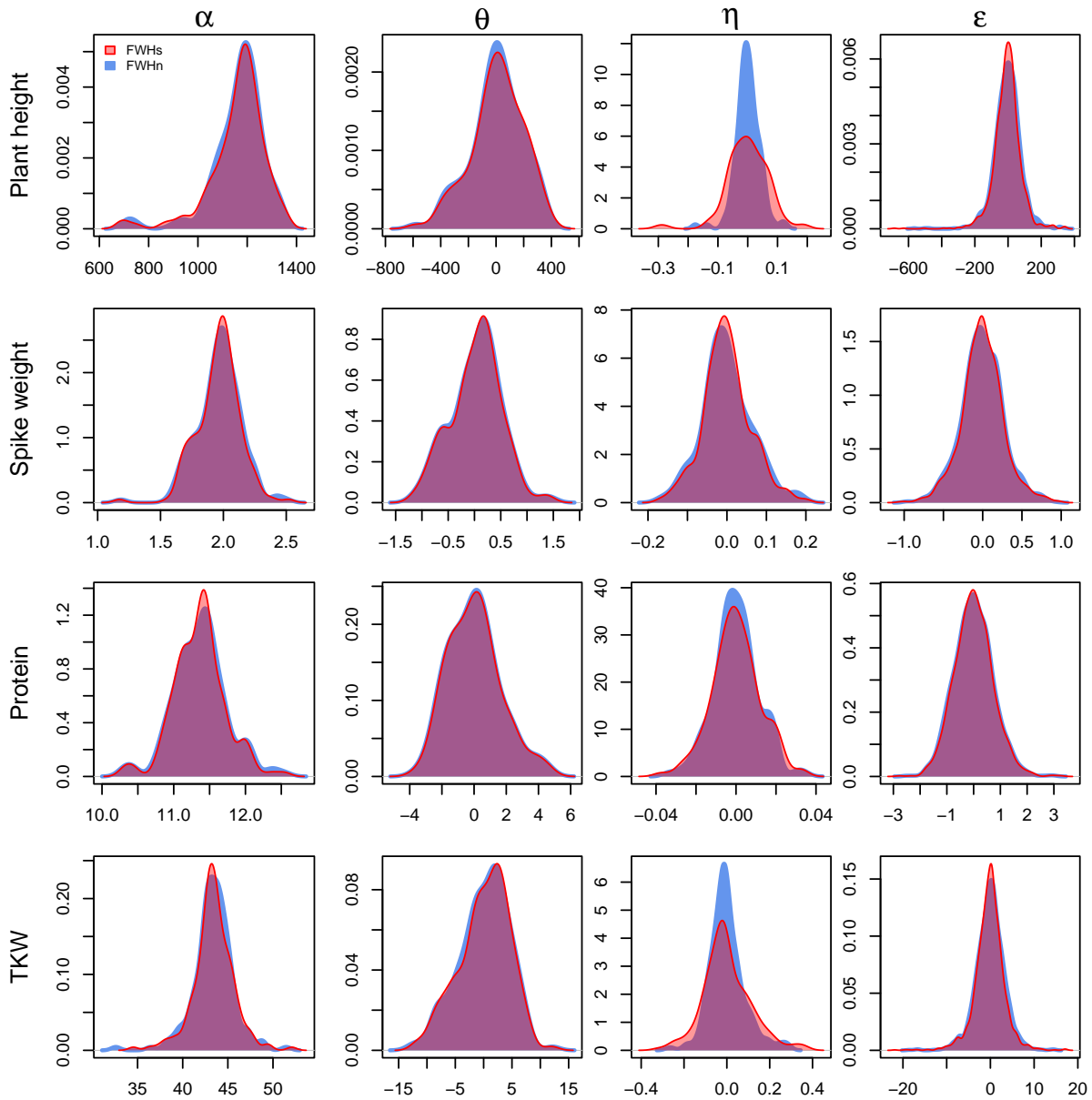


**Figure 2** – The first column presents the distribution of the trait to be explained (in grey). The last four columns compare the hierarchical (red) and non hierarchical (blue) versions of the FW model with a Student law for the residuals, and show the smoothed histograms of main effects, FW coefficients and residuals.

The hierarchical models with a  $t$  distribution (FWHs, ADHs) were more predictive than the models with a normal distribution (FWHn, ADHn), all the more as  $\nu$  was low (Tab. 5). For protein content, the estimate of  $\nu$  was equal to 20, so the  $t$  distribution was close to a normal distribution. The  $t$  distribution reduced the shrinkage of FW coefficients (Fig. 3). Moreover,  $t$  models better accounted for extreme data than normal models (Fig. 4). These extreme data mainly came from germplasm that were not replicated in the trials.

The Finlay-Wilkinson models (FWHs, FWHn) were slightly more predictive than the simple additive models (ADHs, ADHn), except for protein content, where the difference was not significant (Fig. 1). This difference was smaller than the differences due to the distribution of residuals and the hierarchization of parameters.

The  $\text{elpd}_{100}$  criterion was estimated using Pareto smoothed importance sampling (Vehtari et al., 2017). This method tends to be less precise for models that do not fit the data well. Thus, as expected, estimates of  $\text{elpd}_{100}$  were more reliable for the two hierarchical models with a  $t$

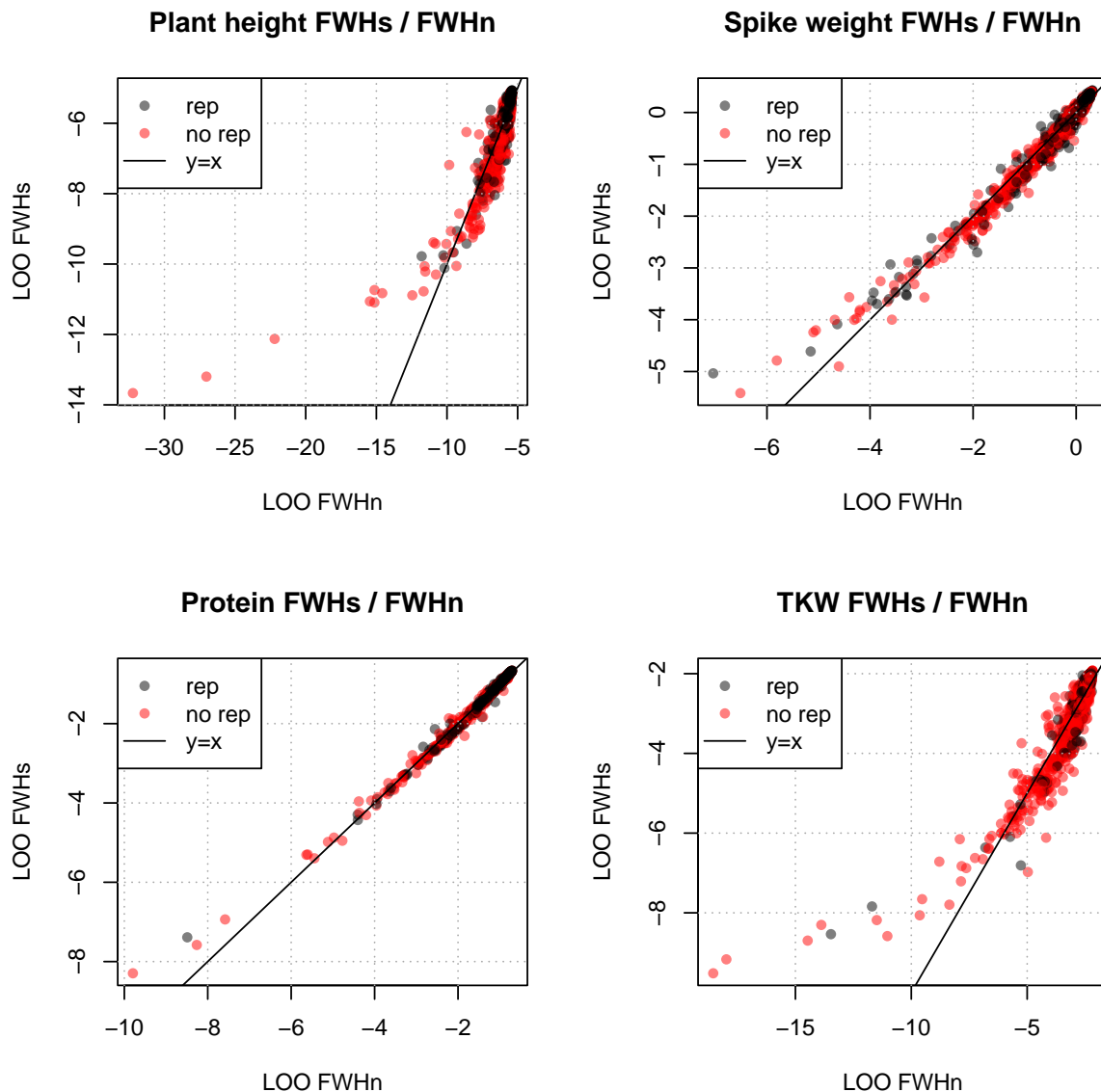


**Figure 3** – Comparison of hierarchical FW models with different residual laws, the Student (red) and the normal (blue). These graphics show the smoothed histograms of main effects, FW coefficients and residuals.

likelihood (FWHs and ADHs) than for the other models, in particular model FWs (Supplementary Tab. B.1).

### 3.2. Precision of estimates and distribution of residuals

For the models with a  $t$  distribution, the estimate of the number of degrees of freedom ( $\nu$ ) varied between 3.4 and 27.6 (close to a normal distribution) (Tab. 5). Thus, the shape of the distribution of residuals depended on the trait. This result confirmed that the number of extreme observations was not negligible in our data, and that models with a  $t$  distribution were more appropriate. In the latter case, the variation ranges of residuals were wider but with more residual values close to 0 for the  $t$  distribution than the normal distribution (Fig. 3). Models had similar estimated precision, except for model FWs, which had less precise estimates. This result confirmed that a basic joint regression, i.e. non-hierarchical model, was not suited to our unbalanced data. Parameters  $\alpha$  and  $\theta$  were estimated more precisely (difference in coefficient of variation between 0 and 1.9, Tab. 5) for  $t$  models (ADHs and FWHs) than for normal models (ADHn and FWHn).



**Figure 4** – Comparison of  $t$  and normal models (FWHs vs FWHn) in terms of the contributions of observations to the  $\text{elpd}_{\text{loo}}$  criterion. Black (resp. red) dots correspond to observations that were measured on germplasm that were replicated (resp. not replicated) within trials.

This result was consistent with Fig. 4, where extreme observations were better predicted by model FWHs than by FWHn, except for protein content.

### 3.3. Variance decomposition

The proportion of variance explained by each term of model FWHs depended on the trait (Tab. 6). For all four traits, the environment effect was highly explanatory. For height and TKW, a relatively large part of the total variance was explained by the germplasm effect (resp. 24% and 16.1%), whereas this part was much smaller for spike weight and protein content (10.9% and 5.7%). The proportion of variance explained by the sensitivity effect  $\eta$  was not significantly different from 0 for protein content and low for the three other traits. It explained 6.7%, 4.9% and 6.9% of the variance of  $G \times E$  interactions and experimental errors ( $\rho$  parameter) for plant height, spike weight and TKW, respectively.

**Table 5** – Number of degrees of freedom and precision of estimates.  $\nu$ : posterior means, with posterior standard deviations in parentheses, of the number of degrees of freedom of the  $t$  distribution;  $cv(\alpha)$ ,  $cv(\theta)$ : average posterior coefficients of variation of germplasm and environment main effects;  $sd(\eta)$ : average posterior standard deviation of germplasm sensitivities (FW coefficients).

Trait	Model	$\nu$	$cv(\alpha)$	$cv(\theta)$	$sd(\eta)$
Plant Height	ADHn		5	4.9	
	ADHs	3.9 (0.5)	3.1	3	
	FWHn		3	2.9	0.08
	FWHs	3.5 (0.4)	2.8	2.7	0.09
	FWs	3.4 (0.4)	3.4	2.7	0.23
Spike weight	ADHn		5.3	5.2	
	ADHs	8.2 (2.3)	5.2	5.1	
	FWHn		5.4	5.2	0.12
	FWHs	8.2 (2.3)	5.2	5.1	0.11
	FWs	10.2 (4)	6.3	4.8	0.31
Protein	ADHn		2.7	2.7	
	ADHs	20.3 (9.8)	2.6	2.7	
	FWHn		2.6	2.7	0.05
	FWHs	19.8 (9.5)	2.6	2.6	0.05
	FWs	27.6 (13)	3.7	2.6	0.27
TKW	ADHn		2.8	2.8	
	ADHs	4.2 (0.5)	2.7	2.5	
	FWHn		2.8	2.8	0.15
	FWHs	4.1 (0.5)	2.7	2.5	0.17
	FWs	3.9 (0.5)	3.1	2.5	0.35

**Table 6** – Variance decomposition for model FWHs. The proportions of variance explained are expressed in %. Mean: posterior mean; 95% CI: 95% credible intervals.  $R^2$  is the coefficient of determination.  $\pi(\alpha)$ ,  $\pi(\theta)$  and  $\pi(\eta\theta)$  are respectively the proportions of variance explained by  $\alpha$ ,  $\theta$  and  $\eta\theta$ .  $\rho$  is the proportion of the variance of  $G \times E$  and errors explained by  $\eta\theta$ .

	Plant height		Spike weight		Protein		TKW	
	Mean	95% CI	Mean	95% CI	Mean	95% CI	Mean	95% CI
$R^2$	87.3	[83.3, 90.3]	78.1	[73.9, 81.9]	82.9	[78.9, 86.6]	69.8	[64.2, 74.9]
$\pi(\alpha)$	24	[18, 30.8]	10.9	[7.8, 14.6]	5.7	[3.7, 8.3]	16.1	[12.1, 20.8]
$\pi(\theta)$	62.4	[54.6, 69.8]	66	[60.1, 71.6]	77	[71.7, 81.9]	51.5	[44.7, 58.2]
$\pi(\eta\theta)$	0.9	[0.4, 1.6]	1.1	[0.4, 2.1]	0.2	[0, 0.8]	2.2	[1, 3.8]
$\rho$	6.7	[2.9, 12.1]	4.9	[1.7, 9.3]	1.2	[0, 4.6]	6.9	[3.1, 12]

### 3.4. Characterization of germplasm

The correlation between germplasm sensitivity ( $\eta_i$ ) and static stability ( $S_i^2$ ) was very close to 1 for all traits while germplasm sensitivity was poorly correlated to  $W_i$  (Tab. 7). The main effect  $\alpha_i$  had a low correlation with  $\eta_i$ ,  $S_i^2$  and  $W_i$ , except for plant height and in some cases spike weight. Depending on the trait, the correlations between  $W_i$  and  $\eta_i$  or  $S_i^2$  were either positive, negative or not significant.

Plant height was found to depend on the type of germplasm, landraces being taller than historic varieties, which were themselves taller than registered varieties. For this trait, registered varieties were significantly more stable (static stability and FW coefficient) than landraces and varieties from crosses, but less stable dynamically (ecoalalance). In addition, registered varieties had lower protein content than the other germplasm types. Landraces and varieties from crosses had lower spike weight than the other germplasm types. Finally, landraces had lower TKW, and historical varieties were statically less stable than the other germplasm types.

**Table 7** – Correlation between germplasm parameters. \*, \*\*, \*\*\* : significant at  $P = 0.05, P = 0.01, P = 0.001$  respectively.  $\alpha_i$ : germplasm effect,  $\eta_i$ : germplasm sensitivity (FW coefficient),  $S_i^2$ : static stability,  $W_i$ : ecovalence.

Trait	Pearson correlation between					
	$\alpha_i \eta_i$	$\alpha_i S_i^2$	$\alpha_i W_i$	$\eta_i S_i^2$	$\eta_i W_i$	$S_i^2 W_i$
Plant height	0.44***	0.41***	-0.43***	0.997***	-0.31***	-0.23**
Spike weight	0.35***	0.35***	0.21**	0.999***	0.15*	0.2**
Protein	0.14	0.13	-0.09	1***	0.03	0.04
TKW	0.23**	0.24**	0.13	0.995***	0.24***	0.34***

**Table 8** – Performance and stability of types of germplasm. For a given line, types with the same letter are not significantly different.  $\bar{\alpha}_k$ : mean germplasm effect of type  $k$ ,  $\bar{\eta}_k$ : mean sensitivity (FW coefficient) of type  $k$ ,  $\bar{S}_k^2$ : mean static stability of type  $k$ , and  $\bar{W}_k$ : mean ecovalence of type  $k$ . This table gives the posterior mean and the 95% credible interval of each parameter.

Trait	Registered	Historic	Landrace	Cross	Mixture	
Plant height	$\bar{\alpha}_k$	862 <sup>d</sup> [822, 901]	1136 <sup>c</sup> [1096, 1175]	1220 <sup>a</sup> [1181, 1258]	1175 <sup>b</sup> [1138, 1210]	1188 <sup>b</sup> [1147, 1228]
	$\bar{\eta}_k$	-0.11 <sup>b</sup> [-0.2, -0.03]	-0.01 <sup>a</sup> [-0.06, 0.05]	0 <sup>a</sup> [-0.05, 0.05]	0.01 <sup>a</sup> [-0.01, 0.04]	-0.01 <sup>a</sup> [-0.09, 0.06]
	$\bar{S}_k^2$	37688 <sup>b</sup> [28956, 48901]	44351 <sup>ab</sup> [34878, 56497]	44813 <sup>a</sup> [35456, 57056]	45620 <sup>a</sup> [36512, 57472]	43800 <sup>ab</sup> [34101, 56150]
	$\bar{W}_k$	9067 <sup>a</sup> [7272, 11725]	7959 <sup>b</sup> [6568, 10184]	7988 <sup>b</sup> [6593, 10234]	7880 <sup>b</sup> [6541, 10056]	7772 <sup>b</sup> [6419, 10000]
	Spike weight	$\bar{\alpha}_k$	2.02 <sup>b</sup> [1.92, 2.12]	1.98 <sup>ab</sup> [1.89, 2.08]	1.93 <sup>a</sup> [1.85, 2.02]	1.96 <sup>a</sup> [1.87, 2.04]
$\bar{\eta}_k$		0.02 <sup>a</sup> [-0.06, 0.1]	0.02 <sup>a</sup> [-0.03, 0.08]	-0.01 <sup>a</sup> [-0.05, 0.03]	0 <sup>a</sup> [-0.02, 0.03]	-0.01 <sup>a</sup> [-0.07, 0.05]
$\bar{S}_k^2$		0.34 <sup>a</sup> [0.27, 0.42]	0.34 <sup>a</sup> [0.28, 0.41]	0.32 <sup>a</sup> [0.27, 0.39]	0.33 <sup>a</sup> [0.28, 0.39]	0.32 <sup>a</sup> [0.26, 0.39]
$\bar{W}_k$		0.08 <sup>a</sup> [0.08, 0.09]	0.08 <sup>a</sup> [0.08, 0.09]	0.08 <sup>a</sup> [0.08, 0.09]	0.08 <sup>a</sup> [0.08, 0.09]	0.08 <sup>a</sup> [0.08, 0.09]
Protein		$\bar{\alpha}_k$	11.08 <sup>a</sup> [10.73, 11.42]	11.37 <sup>b</sup> [11.06, 11.7]	11.4 <sup>b</sup> [11.09, 11.72]	11.35 <sup>b</sup> [11.04, 11.66]
	$\bar{\eta}_k$	0 <sup>a</sup> [-0.03, 0.04]	0 <sup>a</sup> [-0.02, 0.02]	0 <sup>a</sup> [-0.02, 0.02]	0 <sup>a</sup> [-0.01, 0.01]	0 <sup>a</sup> [-0.03, 0.03]
	$\bar{S}_k^2$	3.4 <sup>a</sup> [2.72, 4.29]	3.4 <sup>a</sup> [2.73, 4.27]	3.39 <sup>a</sup> [2.72, 4.26]	3.4 <sup>a</sup> [2.73, 4.27]	3.4 <sup>a</sup> [2.72, 4.28]
	$\bar{W}_k$	0.62 <sup>a</sup> [0.56, 0.68]	0.61 <sup>a</sup> [0.56, 0.68]	0.62 <sup>a</sup> [0.56, 0.68]	0.62 <sup>a</sup> [0.56, 0.68]	0.61 <sup>a</sup> [0.56, 0.68]
	TKW	$\bar{\alpha}_k$	43.6 <sup>ab</sup> [42.6, 44.6]	43.8 <sup>a</sup> [42.9, 44.8]	43.1 <sup>b</sup> [42.3, 43.9]	43.4 <sup>ab</sup> [42.6, 44.1]
$\bar{\eta}_k$		-0.03 <sup>ab</sup> [-0.14, 0.09]	0.08 <sup>a</sup> [0, 0.17]	-0.03 <sup>b</sup> [-0.09, 0.03]	0.01 <sup>ab</sup> [-0.03, 0.05]	-0.05 <sup>b</sup> [-0.14, 0.04]
$\bar{S}_k^2$		33.4 <sup>ab</sup> [26.9, 41.4]	37.8 <sup>a</sup> [31.6, 45.7]	33.3 <sup>b</sup> [28.2, 39.7]	34.9 <sup>ab</sup> [29.7, 41.2]	32.2 <sup>b</sup> [26.6, 39]
$\bar{W}_k$		13.2 <sup>a</sup> [11.5, 15.6]	13.4 <sup>a</sup> [11.6, 15.7]	13.4 <sup>a</sup> [11.7, 15.7]	13.2 <sup>a</sup> [11.6, 15.5]	13.1 <sup>a</sup> [11.4, 15.4]

### 4. Discussion

To fit the characteristics of PPB trials, i.e., few inter-farm replicates and possible extreme data, we developed several models and we found that the hierarchical Finlay-Wilkinson model with  $t$  residuals was the best for prediction and parameter precision. Then we compared the performance and stability of different germplasm types.

#### 4.1. Handling the data from a highly unbalanced series of trials

As the farmers of the program chose the germplasm they assessed, the data obtained from the series of trials were very unbalanced, with more than 90% of the  $G \times E$  combinations missing. This made the estimation of germplasm main effects and sensitivities difficult. Although the Finlay-Wilkinson model was parsimonious, a basic joint regression with weakly-informative prior distributions (model FWs) was not able to cope with this level of disequilibrium. According to the  $\text{elpd}_{\text{loo}}$  criterion, model FWs was not the best model (Fig. 1). In addition, its estimates had poor precision and it led to extreme sensitivity estimates, with values close to 1 or -1 (Fig. 2). In contrast, hierarchical joint regression appeared more suited to our data structure. Model FWs had the largest  $\text{elpd}_{\text{loo}}$  values for three traits out of four. Placing a hierarchical distribution on sensitivities constrained estimates and brought them closer to 0. This led to more satisfactory sensitivity estimates, since they were well below 1 in absolute value.

Three strategies have previously been used to manage incomplete  $G \times E$  data: i) subset the total dataset to obtain an almost balanced subset for the analysis (Ceccarelli and Grando, 2007), ii) predict missing data with a more or less complex model and use these predictions in the analysis (Kumar et al., 2012; Woyann et al., 2017), and iii) use a model more robust to unbalanced data, provided it complies with model validation conditions (Assis et al., 2018; van Frank et al., 2019). We used the last strategy to maximise the amount of information from the data (less data excluded than in the first strategy) with a one-step process (unlike the second strategy).

Cotes et al. (2006) used a Bayesian approach to estimate FW coefficients in a MET study in order to take prior information on germplasm coming from other studies into account. A similar approach was used by Couto et al. (2015), Foucteau and Denis (2001), and Nascimento et al. (2020) and was found to greatly improve the results. Here, we used little prior information. But in the future, previous evaluation studies may provide stronger prior information on germplasm behaviour.

#### 4.2. Extreme observations

Extreme observations were more frequent in our dataset than expected under the normal distribution for three traits out of four (Fig. 4). For these traits, using a  $t$  distribution increased  $\text{elpd}_{\text{loo}}$  values, and the estimate of the number of degrees of freedom of this distribution was smaller than 10 (Tab. 5). In our application, observations were germplasm means resulting from within-trial analyses rather than plot measurements. Extreme observations could occur for several reasons, for example because of the heterogeneity of within-trial residual variances and replications, because cultivation environments were less controlled, or because a non-negligible part of  $G \times E$  interactions was not captured by the multiplicative term of the FW model. The normal distribution was appropriate for protein content. It is difficult to explain why this trait had fewer extreme observations. A possible explanation could be that the measurement of protein content is more standardized than other trait measurements. For plant height, extreme values occurred only for non-replicated micro-plots with a global measurement and never with data from the average of 25 plants (Sect. 2.2.3), suggesting that the plot measurement is less accurate. For TKW, the kernel count could be affected by broken kernels due to over-drying or incorrect threshing settings leading to an overestimation of the number of kernels in the sample. Another possible explanation is that protein content is less variable under different conditions than plant height and spike weight (Kazakou et al., 2014).

Using a  $t$  distribution did not affect the estimates of germplasm and environment main effects. On the contrary, it improved the estimates of sensitivities. It reduced their shrinkage and allowed the multiplicative term of the FW model to better capture  $G \times E$  interactions (Fig. 3).

The Student distribution is expected to take better account of extreme data and to yield more robust estimates (Besag and Higdon, 1999; Lange et al., 1989; Rosa et al., 2003). Extreme data are more likely to occur when varieties are not replicated within trials, which is frequent in this dataset (Fig. 4). Rosa et al. (2003) found that a normal likelihood misestimated a main effect compared to a  $t$  likelihood. This effect was estimated less precisely with a normal distribution, which is consistent with our results for plant height, spike weight and TKW. A Student

distribution appears to be a good solution for dealing with extreme data, in particular in stability analyses, where extreme observations are sometimes removed (this is justified when they are extreme because of experimental errors, but not when they are due to natural variability). While this distribution has recently been used to implement robust alternatives to BLUP (Gianola et al., 2018) or to handle environmental heterogeneity in a single trial (Cao et al., 2022), to our knowledge, it has not already been used in MET studies.

### 4.3. Computing time

Series of trials often include many genotypes and environments, leading to large data sets. Thus, their analysis using mixed or hierarchical models is generally computationally demanding (Smith et al., 2005). The computational load can be reduced by using approximate estimation methods (Nabugoomu et al., 1999) or efficient algorithms, such as algorithms based on sparse matrix operations (Gilmour et al., 1995; Thompson et al., 2003). Hierarchical joint regression has already been implemented using Gibbs sampling or Jags (Lian and Campos, 2016; van Frank et al., 2020). Our implementation based on Hamiltonian Monte Carlo and Stan was more efficient since it required fewer iterations. It allowed us to analyze large datasets in about 6 minutes.

To reduce computing time, the analyses were carried out in two steps. This two-stage approach analyzed  $G \times E$  means without taking account of their standard error, which can reduce the efficiency of the analysis (Welham et al., 2010; Yates and Cochran, 1938). It would be interesting to develop a one-stage method for analyzing plot measurements, in order to better take account of the heterogeneity of the within-trial residual variances and replications (Rivière et al., 2015a).

### 4.4. Variance decomposition

This article shows how to decompose the variance of observations for hierarchical FW models, and how to define the proportions of variance explained by model terms and the coefficient of determination ( $R^2$ ). These quantities are considered as unknown parameters, which are then estimated from the data (Gelman et al., 2019; Helland, 1987). The coefficient of determination is usually defined as the proportion of the sum of squares accounted for by the model, but  $R^2$  defined in this way may be larger than one in a Bayesian framework (Gelman et al., 2019). Our interpretation of  $R^2$  ensures that its estimate is smaller than one.

This variance decomposition is useful to identify the model terms which are the most important. In our application, the environment effects were the most important, explaining from 51% for TKW to 77% for protein content of the variance of observations (Tab. 6). This result is consistent with the diversity of the cropping environments encountered (soil, climate, cropping practices...) and with previous studies (Lian and Campos, 2016; Patterson and Silvey, 1980; Talbot, 1984).

### 4.5. Germplasm main effects and stabilities

Heritability was significant with plant height > TKW > spike weight > protein. Rivière et al. (2015b) found (with data included in our study) a similar ranking in heritability: plant height > TKW = protein > spike weight. Plant height is known to be quite heritable due to a relatively simple genetic architecture with a few major genes, such as the well known Green Revolution Rht1 and Rht2 genes (Peng et al., 1999). In our study, the presence of both recently registered varieties and varieties dating from before the second World War, very likely led to varieties containing different alleles for these loci and increased variability for height. The decrease in plant height from landraces to historic varieties and registered varieties appears very clearly (Tab. 8) as also found in several studies (Bektas et al., 2016; Cantarel et al., 2021).

FW coefficients explained a low proportion of the total variance (between 0.2% and 2.2%) and a low proportion of the variance of  $G \times E$  interactions and errors (between 1.2% and 6.9%, Tab. 6). We can presume that the explanation of the interaction by the FW parameter is weaker the greater the number of environments, for example 29% with less than 10 environments (12 studies), and 12% with more than 10 environments (11 studies, Brancourt-Hulmel et al., 1997). Other classical models, such as AMMI (additive main effect and multiplicative interaction) or



GGE ( $G + G \times E$ ) models, might explain a larger part of  $G \times E$  interactions. Missing data estimation methods allow these models to be used when the data are highly unbalanced, with up to from 40% unbalanced data for a MET with less than 20 environments to 60% unbalanced data for MET with at least 40 environments (Woyann et al., 2017; Yan, 2013). However, these datasets are more balanced than ours, and, as found by Rodrigues et al. (2011), FW is more robust than AMMI when the data are highly unbalanced (75%). In our study, most germplasm occurred in a limited number of environments, so that a parsimonious and very simple modelling of  $G \times E$  interactions had to be used. An alternative approach would be to better characterize the environments and thus explain the environmental effects and part of the  $G \times E$  interaction using environmental variables (Piepho and Blancon, 2023).

Although sensitivities explained a rather low proportion of variance, FWs model had larger  $elp_{d_{100}}$  values than additive models for three traits out of four. In addition, for these traits, some sensitivity estimates were not negligible, with values close to 0.2 or 0.3. Interaction effects then represented 20% or 30% of environmental effects. Additive models were appropriate for the protein content trait. It was found that the multiplicative term of the FW model was not significant for protein content, both in a balanced network of 15 environments in Serbia (Hristov et al., 2010) and in 12 environments in Swiss organic trials (Knapp et al., 2017). On the contrary, Mut et al. (2010) found significant FW coefficients for a balanced network of 7 environments in Turkey. These contrasting results could be explained by differences between numbers of environments or between genetic diversities.

For plant height, we found that registered varieties were more statically stable but less dynamically stable (Tab. 8). This can be explained by the fact that there are only a few registered varieties in the trials, therefore they have little influence on the average height, which can fluctuate greatly between trials, and therefore the deviation from this average will be greater for this type.

Static and dynamic stabilities were difficult to estimate since our series of trials was very unbalanced. In particular, raw estimates of these stabilities were not reliable, since they were very influenced by the unbalanced nature of the data. By using theoretical variances, the FW model allowed us to calculate simple indicators of static and dynamic stability in the wheat PPB dataset. However, comparisons between germplasm stability indicators only take account of the part of  $G \times E$  interactions explained by the FW model. To our knowledge, the FW model has never been used for this purpose before.

Dependence between stability and mean is widespread (Reckling et al., 2021), but in our case, the correlation was low, which simplified interpretation of the stability analysis. Several studies for different traits and with balanced MET found a very strong correlation between FW coefficient and the static stability (Becker, 1981; Fasahat et al., 2015; Reckling et al., 2021). However, in our case, this relationship was even stronger (Tab. 7), probably because of the assumption that the variance of residuals did not depend on the genotype. As in many other studies, the residual variance was assumed to be independent of germplasm throughout our study. Allowing the residual variance to depend on the genotype could improve the estimates of stability indicators (Cotes et al., 2006; Couto et al., 2015). In particular, the dynamic indicator would be similar to the Shukla Stability Variance, i.e, the varietal variance of  $G \times E$  interactions (Cotes et al., 2006). However, estimating a residual variance and a FW coefficient for each germplasm could be difficult in our study, as most of the germplasm appeared in only a few environments.

When relationships were significant, mixtures were always in a more stable (statically and dynamically) statistical group (Tab. 8). This result supports the fact that within-plot diversity stabilizes performances (Döring et al., 2015; Kiær et al., 2012).

In the wheat PPB program, the populations tested were heterogeneous and their genetic composition could vary over years and farms (David et al., 2020). In this analysis, such variations were considered as part of the response of a population to a given environment for the sake of simplicity. Therefore the  $G \times E$  interactions could be overestimated (resp. underestimated) if populations underwent diversifying (resp. stabilizing) selection pressures within farms.

One aim of the project was to provide farmers with information to help them select new germplasm for testing in their farm. The statistical tools we developed sought to cope with the large degree to which this series of trials was unbalanced. Their objectives were the same as in other MET analyses : (i) estimate and predict germplasm values for traits of interest for breeding, (ii) study the stability of germplasm over several environments, (iii) select new germplasm to be tested in new locations (Cotes et al., 2006). MET are usually carried out to find stable germplasm that perform well on average over many locations, or to detect special local adaptations to certain environments (Annicchiarico et al., 2005; Gauch et al., 2008). Here, while farmers were mostly interested in selecting the best germplasm adapted to their local pedo-climatic conditions, farming practices and marketing objectives, information retrieved from the farmers' network on new varieties to introduce in their trials could also be useful.

## 5. Conclusion

The proposed hierarchical model aims to improve the estimates of the parameters of the FW model from unbalanced datasets. This model was complex and was easier to implement in a Bayesian framework. Placing hierarchical distributions on model parameters and modelling residuals using a  $t$  distribution improved the estimates of main and interaction effects. This model allowed us to estimate static and dynamic stability indicators despite the high level of data imbalance. Main effects and stability indicators provide information on the behaviour of genotypes in different environments, which farmers could use in their selection process.

Participatory research raises new research questions and contributes to the development of new methods for societal action (Kastenhofer et al., 2011). In PPB programs, all the methodology is based on collective and collaborative work and action between farmers, associations of farmers and researchers (Brac de la Perrière et al., 2011). New statistical methods can contribute to a better use of such complex multi-environment data in the selection process, and more generally to the effectiveness of participatory research (Martin and Sherington, 1997).

## Acknowledgements

We thank Estelle Serpolay, Nathalie Galic and Sophie Pin for their great help in the measurements on participating farms and research stations. We thank all the farmers and facilitators participating in the project. We thank Gérard Branlard from INRAE Clermont-Ferrand for his time when carrying out NIRS analysis and Jean-Marc Le Goff for his feedback. We thank Pierre Druilhet and David Makowski for helpful comments.

Preprint version 2 of this article has been peer-reviewed and recommended by Peer Community In Mathematical and Computational Biology (<https://doi.org/10.24072/pci.mcb.100272>); Donnet, 2024).

## Funding

M. Turbet Delof was funded by Program PPR-CPA MoBiDiv (2021–2026) under grant agreement ANR-20-PCPA-0006 and INRAE (program on organic scaling METABIO and *Biologie et Amélioration des Plantes* department).

## Conflict of interest disclosure

The authors declare that they comply with the PCI rule of having no financial conflicts of interest in relation to the content of the article.

## Data, script, code, and supplementary information availability

Data and script for models (version 1) are available online in the server *Recherche Data Gouv* (<https://doi.org/10.57745/SUTZ9U>; Turbet Delof et al., 2024).

### Supplementary information

#### Appendix A. Models

Tab. A.1 provides supplementary information on the prior distribution of model parameters.

**Table A.1** – Known values of the parameters of the prior distribution ( $\lambda_\mu, \lambda_\varepsilon$ ), empirical mean ( $\mu_{emp}$ ) and standard deviation ( $\sigma_{emp}$ ) of traits.

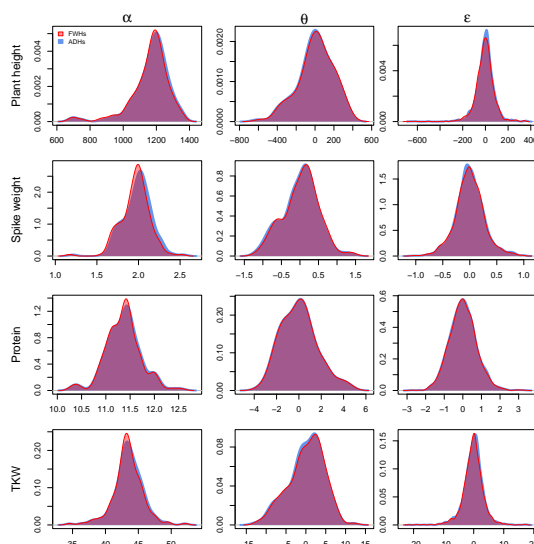
	$\lambda_\mu$	$\lambda_\varepsilon$	$\mu_{emp}$	$\sigma_{emp}$
Plant height	1200	500	1188	234
Spike weight	2.00	0.80	2.03	0.58
Protein	12.0	4.0	11.5	1.9
TKW	45.0	10.0	43.7	5.8

#### Appendix B. Model comparison

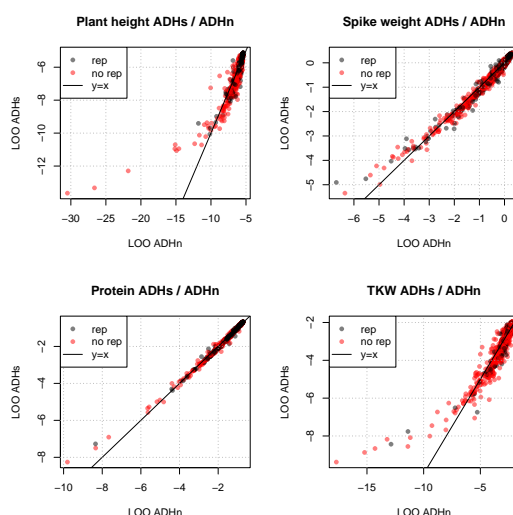
Tab. B.1 provides supplementary information on the estimation of the  $elpd_{100}$  criterion. Fig. B.1 provides supplementary information on the comparison of models FWHs and ADHs. Fig. B.2 provides supplementary information on the comparison of models ADHs and ADHn.

**Table B.1** – Estimates of tail shape parameters ( $k$ ) used to estimate  $elpd_{100}$ . The contribution of each observation to  $elpd_{100}$ , i.e.,  $\ln(p(Y_{ij}|Y_{-ij}))$ , was estimated using Pareto smoothed importance sampling (Vehtari et al., 2017). For each observation, the largest importance weights of the importance sampling were smoothed using a generalized Pareto distribution with shape parameter  $k$ . Estimates of pointwise contributions with  $k > 0.7$  are less reliable.

Trait	Model	$k < 0.5$	$0.5 < k < 0.7$	$0.7 < k < 1$	$k > 1$
Spike weight	ADHn	1396	32	7	2
	ADHs	1437	0	0	0
	FWHn	1397	31	7	2
	FWHs	1437	0	0	0
	FWs	1404	25	7	1
Plant height	ADHn	1781	23	0	0
	ADHs	1804	0	0	0
	FWHn	1761	39	4	0
	FWHs	1804	0	0	0
	FWs	1664	127	12	1
TKW	ADHn	1314	16	2	0
	ADHs	1331	1	0	0
	FWHn	1300	31	1	0
	FWHs	1329	3	0	0
	FWs	1121	150	53	8
Protein	ADHn	1955	26	1	0
	ADHs	1982	0	0	0
	FWHn	1907	69	6	0
	FWHs	1981	1	0	0
	FWs	1937	43	2	0



**Figure B.1** – Comparison of models FWHs and ADHs for the distribution of germplasm main effects ( $\alpha$ ), environment main effects ( $\theta$ ) and residuals ( $\varepsilon$ ) for each trait. Red: model FWHs; Blue: model ADHs.



**Figure B.2** – Comparison of models ADHs and ADHn in terms of the contributions of observations to the  $elpd_{loo}$  criterion. Black (resp. red) dots correspond to observations that were measured on germplasm that were repeated (resp. not repeated) within trials.

## References

- Annicchiarico P (2002). *Genotype x environment interaction: challenges and opportunities for plant breeding and cultivar recommendations*. FAO plant production and protection paper 174. Rome: Food and Agriculture Organization of the United Nations.
- Annicchiarico P, Bellah F, Chiari T (2005). *Defining Subregions and Estimating Benefits for a Specific-Adaptation Strategy by Breeding Programs: A Case Study*. *Crop Science* **45**, 1741–1749. <https://doi.org/10.2135/cropsci2004.0524>.
- Annicchiarico P, Chiapparino E, Perenzin M (2010). *Response of Common Wheat Varieties to Organic and Conventional Production Systems across Italian Locations, and Implications for Selection*. *Field Crops Research* **116**, 230–238. <https://doi.org/10.1016/j.fcr.2009.12.012>.

- Assis TOG, Dias CTDS, Rodrigues PC (2018). A weighted AMMI algorithm for nonreplicated data. *Pesquisa Agropecuária Brasileira* **53**, 557–565. <https://doi.org/10.1590/s0100-204x2018000500004>. URL: [http://www.scielo.br/scielo.php?script=sci\\_arttext&pid=S0100-204X20180005000557&lng=en&tlng=en](http://www.scielo.br/scielo.php?script=sci_arttext&pid=S0100-204X20180005000557&lng=en&tlng=en).
- Becker HC (1981). *Correlations among Some Statistical Measures of Phenotypic Stability*. *Euphytica* **30**, 835–840. <https://doi.org/10.1007/bf00038812>.
- Becker HC, Leon J (1988). *Stability Analysis in Plant Breeding*. *Plant breeding* **101**, 1–23. <https://doi.org/10.1111/j.1439-0523.1988.tb00261.x>.
- Bektas H, Hohn CE, Waines JG (2016). *Root and shoot traits of bread wheat (Triticum aestivum L.) landraces and cultivars*. *Euphytica* **212**, 297–311. <https://doi.org/10.1007/s10681-016-1770-7>. URL: <https://link.springer.com/10.1007/s10681-016-1770-7>.
- Besag J, Higdon D (1999). *Bayesian Analysis of Agricultural Field Experiments*. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **61**, 691–746. <https://doi.org/10.1111/1467-9868.00201>.
- Beza E, Steinke J, Van Etten J, Reidsma P, Fadda C, Mitra S, Mathur P, Kooistra L (2017). *What are the prospects for citizen science in agriculture? Evidence from three continents on motivation and mobile telephone use of resource-poor farmers*. *PloS one* **12**, e0175700. <https://doi.org/10.1371/journal.pone.0175700>.
- Brac de la Perrière RA, De Kochko P, Neubauer C, Storup B (2011). *Visions Paysannes de La Recherche Dans Le Contexte de La Sélection Participative*. Emergence. Pour l'émergence d'une université du Vivant.
- Brancourt-Hulmel M, Biarnès-Dumoulin V, Denis JB (1997). *Points de repère dans l'analyse de la stabilité et de l'interaction génotype-milieu en amélioration des plantes*. *Agronomie* **17**, 219–246. <https://doi.org/10.1051/agro:19970403>.
- Cantarel AA, Allard V, Andrieu B, Barot S, Enjalbert J, Gervais J, Goldringer I, Pommier T, Saint-Jean S, Le Roux X (2021). *Plant Functional Trait Variability and Trait Syndromes among Wheat Varieties: The Footprint of Artificial Selection*. *Journal of Experimental Botany* **72**, 1166–1180. <https://doi.org/10.1093/jxb/eraa491>.
- Cao Z, Stefanova K, Gibberd M, Rakshit S (2022). *Bayesian inference of spatially correlated random parameters for on-farm experiment*. *Field Crops Research* **281**, 108477. <https://doi.org/10.1016/j.fcr.2022.108477>. URL: <https://linkinghub.elsevier.com/retrieve/pii/S037842902200048X>.
- Carlin BP, Polson NG (1991). *Inference for Nonconjugate Bayesian Models Using the Gibbs Sampler*. *Canadian Journal of Statistics* **19**, 399–405. <https://doi.org/10.2307/3315430>.
- Carlin BP, Louis TA (2009). *Bayesian methods for data analysis*. Third edition. Chapman & Hall/CRC Texts in Statistical Science Series. Boca Raton, Florida ; London, England ; New York, New York: CRC Press.
- Ceccarelli S, Grando S, Bailey E, Amri A, El-Felah M, Nassif F, Rezgui S, Yahyaoui A (2001). *Farmer Participation in Barley Breeding in Syria, Morocco and Tunisia*. *Euphytica* **122**, 521–536. <https://doi.org/10.1023/A:1017570702689>.
- Ceccarelli S, Grando S (2007). *Decentralized-Participatory Plant Breeding: An Example of Demand Driven Research*. *Euphytica* **155**, 349–360. <https://doi.org/10.1007/s10681-006-9336-8>.
- Ceccarelli S, Grando S (2020). *Participatory Plant Breeding: Who Did It, Who Does It and Where?* *Experimental Agriculture* **56**, 1–11. <https://doi.org/10.1017/s0014479719000127>.
- Choy BST, Chan JSK (2008). *Scale Mixtures Distributions in Statistical Modelling*. *Australian & New Zealand Journal of Statistics* **50**, 135–146. <https://doi.org/10.1111/j.1467-842x.2008.00504.x>.
- Cotes JM, Crossa J, Sanches A, Cornelius PL (2006). *A Bayesian Approach for Assessing the Stability of Genotypes*. *Crop Science* **46**, 2654–2665. <https://doi.org/10.2135/cropsci2006.04.0227>.
- Couto MF, Nascimento M, Amaral Jr AT, Silva FF, Viana AP, Vivas M (2015). *Eberhart and Russel's Bayesian Method in the Selection of Popcorn Cultivars*. *Crop Science* **55**, 571–577. <https://doi.org/10.2135/cropsci2014.07.0498>.

- David O, van Frank G, Goldringer I, Rivière P, Turbet Delof M (2020). *Bayesian Inference of Natural Selection from Spatiotemporal Phenotypic Data*. *Theoretical Population Biology* **131**, 100–109. <https://doi.org/10.1016/j.tpb.2019.11.007>.
- Dawson J, Murphy KM, Jones SS (2008). *Decentralized Selection and Participatory Approaches in Plant Breeding for Low-Input Systems*. *Euphytica* **160**, 143–154. <https://doi.org/10.1007/s10681-007-9533-0>.
- Dawson JC, Rivière P, Berthelot JF, Mercier F, Kochko P, Galic N, Pin S, Serpolay E, Thomas M, Giuliano S, Goldringer I (2011). *Collaborative Plant Breeding for Organic Agricultural Systems in Developed Countries*. *Sustainability* **3**, 1206–1223. <https://doi.org/10.3390/su3081206>.
- Desclaux D, Nolot JM, Chiffolleau Y, Gozé E, Leclerc C (2008). *Changes in the Concept of Genotype × Environment Interactions to Fit Agriculture Diversification and Decentralized Participatory Plant Breeding: Pluridisciplinary Point of View*. *Euphytica* **163**, 533–546. <https://doi.org/10.1007/s10681-008-9717-2>.
- Donnet S (2024). *Handling Data Imbalance and G×E Interactions in On-Farm Trials Using Bayesian Hierarchical Models*. *Peer Community in Mathematical and Computational Biology* **1**, 100272. <https://doi.org/10.24072/pci.mcb.100272>.
- Döring TF, Annicchiarico P, Clarke S, Haigh Z, Jones HE, Pearce H, Snape J, Zhan J, Wolfe MS (2015). *Comparative analysis of performance and stability among composite cross populations, variety mixtures and pure lines of winter wheat in organic and conventional cropping systems*. *Field Crops Research* **183**, 235–245. <https://doi.org/10.1016/j.fcr.2015.08.009>. URL: <https://www.sciencedirect.com/science/article/pii/S037842901530040X>.
- Falconer DS (1960). *Introduction to Quantitative Genetics*. Edinburgh/London: Oliver & Boyd.
- Fasahat P, Rajabi A, Mahmoudi SB, Noghabi MA, Rad JM (2015). *An Overview on the Use of Stability Parameters in Plant Breeding*. *Biometrics & Biostatistics International Journal* **2**, 149–159. <https://doi.org/10.15406/bbij.2015.02.00043>.
- Finlay KW, Wilkinson GN (1963). *The Analysis of Adaptation in a Plant-Breeding Programme*. *Australian Journal of Agricultural Research* **14**, 742–754. <https://doi.org/10.1071/AR9630742>.
- Foucteau V, Denis JB (2001). *Statistical Analysis of Successive Series of Experiments in Plant Breeding: A Bayesian Approach*. In: Paris, France: Institut National de la Recherche Agronomique, pp. 49–56.
- Gauch HG, Piepho HP, Annicchiarico P (2008). *Statistical Analysis of Yield Trials by AMMI and GGE: Further Considerations*. *Crop Science* **48**, 866–889. <https://doi.org/10.2135/cropsci2007.09.0513>.
- Gelman A (2006). *Prior distributions for variance parameters in hierarchical models (comment on article by Browne and Draper)*. *Bayesian Analysis* **1**, 515–534. <https://doi.org/10.1214/06-BA117A>. URL: <https://projecteuclid.org/journals/bayesian-analysis/volume-1/issue-3/Prior-distributions-for-variance-parameters-in-hierarchical-models-comment-on/10.1214/06-BA117A.full>.
- Gelman A, Goodrich B, Gabry J, Vehtari A (2019). *R-squared for Bayesian Regression Models*. *The American Statistician* **73**, 307–309. <https://doi.org/10.1080/00031305.2018.1549100>.
- Gianola D, Cecchinato A, Naya H, Schön CC (2018). *Prediction of complex traits: robust alternatives to best linear unbiased prediction*. *Frontiers in genetics* **9**, 195. <https://doi.org/10.3389/fgene.2018.00195>.
- Gilmour AR, Thompson R, Cullis BR (1995). *Average Information REML: An Efficient Algorithm for Variance Parameter Estimation in Linear Mixed Models*. *Biometrics*, 1440–1450. <https://doi.org/10.2307/2533274>.
- Goldringer I, van Frank G, Bouvier d'Yvoire C, Forst E, Galic N, Garnault M, Locqueville J, Pin S, Bailly J, Baltassat R, Berthelot JF, Caizergues F, Dalmaso C, Kochko P, Gascuel JS, Hyacinthe A, Lacanette J, Mercier F, Montaz H, Ronot B, et al. (2020). *Agronomic Evaluation of Bread Wheat Varieties from Participatory Breeding: A Combination of Performance and Robustness*. *Sustainability* **12**, 128. <https://doi.org/10.3390/su12010128>.
- Hampel FR, Ronchetti EM, Rousseeuw PJ, Stahel WA (2011). *Robust Statistics: The Approach Based on Influence Functions*. John Wiley & Sons. <https://doi.org/10.1002/9781118186435>.

- Helland IS (1987). *On the Interpretation and Use of  $R^2$  in Regression Analysis*. *Biometrics* **43**, 61–69. <https://doi.org/10.2307/2531949>. URL: <https://www.jstor.org/stable/2531949?origin=crossref>.
- Hristov N, Mladenov N, Djuric V, Kondic-Spika A, Marjanovic-Jeromela A, Simic D (2010). *Genotype by Environment Interactions in Wheat Quality Breeding Programs in Southeast Europe*. *Euphytica* **174**, 315–324. <https://doi.org/10.1007/s10681-009-0100-8>.
- Huber PJ, Ronchetti EM (1981). *Robust Statistics*. John Wiley & Sons. <https://doi.org/10.1002/9780470434697>.
- Juárez MA, Steel MFJ (2010). *Model-Based Clustering of Non-Gaussian Panel Data Based on Skew-t Distributions*. *Journal of Business & Economic Statistics* **28**, 52–66. <https://doi.org/10.1198/jbes.2009.07145>.
- Kastenhofer K, Bechtold U, Wilfing H (2011). *Sustaining Sustainability Science: The Role of Established Inter-Disciplines*. *Ecological Economics* **70**, 835–843. <https://doi.org/10.1016/j.ecolecon.2010.12.008>.
- Kazakou E, Violle C, Roumet C, Navas ML, Vile D, Kattge J, Garnier E (2014). *Are Trait-based Species Rankings Consistent across Data Sets and Spatial Scales?* *Journal of Vegetation Science* **25**, 235–247. <https://doi.org/10.1111/jvs.12066>.
- Kiær LP, Skovgaard IM, Østergård H (2012). *Effects of Inter-Varietal Diversity, Biotic Stresses and Environmental Productivity on Grain Yield of Spring Barley Variety Mixtures*. *Euphytica* **185**, 123–138. <https://doi.org/10.1007/s10681-012-0640-1>.
- Knapp S, Brabant C, Oberforster M, Grausgruber H, Hiltbrunner J (2017). *Quality Traits in Winter Wheat: Comparison of Stability Parameters and Correlations between Traits Regarding Their Stability*. *Journal of Cereal Science* **77**, 186–193. <https://doi.org/10.1016/j.jcs.2017.08.011>.
- Kumar A, Verulkar SB, Mandal NP, Variar M, Shukla VD, Dwivedi JL, Singh BN, Singh ON, Swain P, Mall AK (2012). *High-Yielding, Drought-Tolerant, Stable Rice Genotypes for the Shallow Rainfed Lowland Drought-Prone Ecosystem*. *Field crops research* **133**, 37–47. <https://doi.org/10.1016/j.fcr.2012.03.007>.
- Lange KL, Little RJ, Taylor JM (1989). *Robust Statistical Modeling Using the t Distribution*. *Journal of the American Statistical Association* **84**, 881–896. <https://doi.org/10.2307/2290063>.
- Lartillot N (2023). *Identifying the Best Approximating Model in Bayesian Phylogenetics: Bayes Factors, Cross-Validation or wAIC?* *Systematic Biology* **72**. Ed. by Sebastian Höhna, 616–638. <https://doi.org/10.1093/sysbio/syad004>. URL: <https://academic.oup.com/sysbio/article/72/3/616/7050004>.
- Lian L, Campos G (2016). *FW: An R Package for Finlay–Wilkinson Regression That Incorporates Genomic/Pedigree Information and Covariance Structures Between Environments*. *G3 Genes|Genomes|Genetics* **6**, 589–597. <https://doi.org/10.1534/g3.115.026328>.
- Lin CS, Binns MR, Lefkovitch LP (1986). *Stability Analysis: Where Do We Stand?* *Crop science* **26**, 894–900. <https://doi.org/10.2135/cropsci1986.0011183x002600050012x>.
- Martin A, Sherington J (1997). *Participatory Research Methods— Implementation, Effectiveness and Institutional Context*. *Agricultural systems* **55**, 195–216. [https://doi.org/10.1016/s0308-521x\(97\)00007-3](https://doi.org/10.1016/s0308-521x(97)00007-3).
- Mohammadi R, Mahmoodi KN, Haghparast R, Grando S, Rahmanian M, Ceccarelli S (2011). *Identifying Superior Rainfed Barley Genotypes in Farmers' Fields Using Participatory Varietal Selection*. *Journal of Crop Science and Biotechnology* **14**, 281–288. <https://doi.org/10.1007/s12892-010-0106-8>.
- Murphy KM, Campbell KG, Lyon SR, Jones SS (2007). *Evidence of Varietal Adaptation to Organic Farming Systems*. *Field Crops Research* **102**, 172–177. <https://doi.org/10.1016/j.fcr.2007.03.011>.
- Mut Z, Aydin N, Orhan Bayramoglu H, Ozcan H (2010). *Stability of Some Quality Traits in Bread Wheat (*Triticum Aestivum*) Genotypes*. *Journal of Environmental Biology* **31**, 489.
- Nabugoomu F, Kempton RA, Talbot M (1999). *Analysis of Series of Trials Where Varieties Differ in Sensitivity to Locations*. *Journal of Agricultural, Biological, and Environmental Statistics* **4**, 310–325. <https://doi.org/10.2307/1400388>.

- Nascimento M, Nascimento ACC, Silva FF, Teodoro PE, Azevedo CF, Oliveira TRA, Amaral Junior AT, Cruz CD, Farias FJC, Carvalho LP (2020). *Bayesian Segmented Regression Model for Adaptability and Stability Evaluation of Cotton Genotypes*. *Euphytica* **216**, 1–10. <https://doi.org/10.1007/s10681-020-2564-5>.
- Neal RM (2011). *MCMC Using Hamiltonian Dynamics*. In: *Handbook of Markov Chain Monte Carlo*. Ed. by Steve Brooks, Andrew Gelman, Galin Jones, and Xiao-Li Meng. Chapman and Hall/CRC, pp. 113–162. <https://doi.org/10.1201/b10905-6>.
- Ng M, Williams E (2001). *Joint-Regression Analysis for Incomplete Two-way Tables*. *Australian & New Zealand Journal of Statistics* **43**, 201–206. <https://doi.org/10.1111/1467-842X.00165>.
- Patterson HD (1997). *Analysis of series of variety trials*. In: *Statistical methods for plant variety evaluation*. Ed. by Rodney Alistair Kempton and Paul N Fox. Chapman & Hall, pp. 139–161. [https://doi.org/10.1007/978-94-009-1503-9\\_9](https://doi.org/10.1007/978-94-009-1503-9_9).
- Patterson HD, Silvey V (1980). *Statutory and Recommended List Trials of Crop Varieties in the United Kingdom*. *Journal of the Royal Statistical Society: Series A (General)* **143**, 219–240. <https://doi.org/10.2307/2982128>.
- Peng J, Richards DE, Hartley NM, Murphy GP, Devos KM, Flintham JE, Beales J, Fish LJ, Worland AJ, Pelica F (1999). 'Green Revolution' Genes Encode Mutant Gibberellin Response Modulators. *Nature* **400**, 256–261. <https://doi.org/10.1038/22307>.
- Pereira DG, Mexia JT, Rodrigues PC (2007). *Robustness of Joint Regression Analysis*. *Biometrical Letters* **44**, 105–128. <https://doi.org/10.1590/s0103-90162011000600012>.
- Perkins JM, Jinks JL (1968). *Environmental and genotype-environmental components of variability III. Multiple lines and crosses*. *Heredity* **23**, 339–356. <https://doi.org/10.1038/hdy.1968.48>. URL: <https://www.nature.com/articles/hdy196848>.
- Piepho HP (1999). *Stability analysis using the SAS system*. *Agronomy Journal* **91**, 154–160. <https://doi.org/10.2134/agronj1999.00021962009100010024x>.
- Piepho HP (2004). *An algorithm for a letter-based representation of all-pairwise comparisons*. *Journal of Computational and Graphical Statistics* **13**, 456–466. <https://doi.org/10.1002/bimj.200490128>.
- Piepho HP, Blancon J (2023). *Extending Finlay–Wilkinson regression with environmental covariates*. *Plant Breeding* **142**, 621–631. <https://doi.org/10.1111/pbr.13130>.
- R Core Team (2024). *R: A Language and Environment for Statistical Computing*. Place: Vienna, Austria. URL: <http://www.R-project.org/>.
- Reckling M, Ahrends H, Chen TW, Eugster W, Hadasch S, Knapp S, Laidig F, Linstädter A, Macholdt J, Piepho HP (2021). *Methods of Yield Stability Analysis in Long-Term Field Experiments. A Review*. *Agronomy for Sustainable Development* **41**, 1–28. <https://doi.org/10.1007/s13593-021-00681-4>.
- Rivière P, Dawson JC, Goldringer I, David O (2015a). *Hierarchical Bayesian Modeling for Flexible Experiments in Decentralized Participatory Plant Breeding*. *Crop Science* **55**, 1053–1067. <https://doi.org/10.2135/cropsci2014.07.0497>.
- Rivière P, Goldringer I, Berthelot JF, Galic N, Pin S, De Kochko P, Dawson JC (2015b). *Response to Farmer Mass Selection in Early Generation Progeny of Bread Wheat Landrace Crosses*. *Renewable Agriculture and Food Systems* **30**, 190–201. <https://doi.org/10.1017/S1742170513000343>.
- Robert C (2007). *The Bayesian Choice: From Decision-Theoretic Foundations to Computational Implementation*. Springer. <https://doi.org/10.1007/0-387-71599-1>.
- Rodrigues PC, Pereira DGS, Mexia JT (2011). *A Comparison between Joint Regression Analysis and the Additive Main and Multiplicative Interaction Model: The Robustness with Increasing Amounts of Missing Data*. *Scientia Agricola* **68**, 679–686. <https://doi.org/10.1590/s0103-90162011000600012>.
- Rosa G, Padovani CR, Gianola D (2003). *Robust Linear Mixed Models with Normal/Independent Distributions and Bayesian MCMC Implementation*. *Biometrical Journal* **45**, 573–590. <https://doi.org/10.1002/bimj.200390034>.
- Silvey SD (1975). *Statistical Inference*. Chapman and Hall.



- Simar L (2002). *Les modèles de base de l'analyse bayésienne*. In: *Méthodes Bayésiennes en statistique*. Ed. by Jean-Jacques Dreesbeke, Jeanne Fine, and Gilbert Saporta. Editions Technip, pp. 103–132. <https://doi.org/10.1787/dev-cred-data-fr>.
- Smith A, Cullis BR, Thompson R (2005). *The Analysis of Crop Cultivar Breeding and Evaluation Trials: An Overview of Current Mixed Model Approaches*. *The Journal of Agricultural Science* **143**, 449–462. <https://doi.org/10.1017/s0021859605005587>.
- Smith ME, Castillo FG, Gómez F (2001). *Participatory Plant Breeding with Maize in Mexico and Honduras*. *Euphytica* **122**, 551–563. <https://doi.org/10.1023/A:1017510529440>.
- Snapp SS, Silim SN (2002). *Farmer Preferences and Legume Intensification for Low Nutrient Environments*. *Plant and soil* **245**, 181–192. [https://doi.org/10.1007/978-94-017-1570-6\\_31](https://doi.org/10.1007/978-94-017-1570-6_31).
- Stan Development Team (2024). *RStan: The R Interface to Stan*. URL: <http://mc-stan.org/>.
- Talbot M (1984). *Yield variability of crop varieties in the UK*. *The Journal of Agricultural Science* **102**, 315–321. <https://doi.org/10.1017/s0021859600042635>.
- Thompson R, Cullis B, Smith A, Gilmour A (2003). *A Sparse Implementation of the Average Information Algorithm for Factor Analytic and Reduced Rank Variance Models*. *Australian & New Zealand Journal of Statistics* **45**, 445–459. <https://doi.org/10.1111/1467-842x.00297>.
- Turbet Delof M (2024). *Impacts de l'environnement sur les pratiques de sélection paysanne et le comportement des variétés qui en résultent*. PhD Thesis. Université Paris-Saclay. <https://doi.org/10.1787/1876f33f-fr>.
- Turbet Delof M, Rivière P, Dawson JC, Gauffreteau A, Goldringer I, van Frank G, David O (2024). *Supplementary Data and models : Bayesian joint-regression analysis of unbalanced series of on-farm trials*. <https://doi.org/10.57745/SUTZ9U>.
- Van Etten J, De Sousa K, Aguilar A, Barrios M, Coto A, Dell'Acqua M, Fadda C, Gebrehawaryat Y, Van De Gevel J, Gupta A, Kiros AY, Madriz B, Mathur P, Mengistu DK, Mercado L, Nurhisen Mohammed J, Paliwal A, Pè ME, Quirós CF, Rosas JC, et al. (2019). *Crop variety management for climate adaptation supported by citizen science*. *Proceedings of the National Academy of Sciences* **116**, 4194–4199. <https://doi.org/10.1073/pnas.1813720116>. URL: <https://pnas.org/doi/full/10.1073/pnas.1813720116>.
- van Frank G (2018). *Gestion participative de la diversité cultivée et création de mélanges diversifiés de blé tendre à la ferme*. PhD Thesis. Université Paris Saclay (COMUE). <https://doi.org/10.1522/12332444>.
- van Frank G, Goldringer I, Rivière P, David O (2019). *Influence of Experimental Design on Decentralized, on-Farm Evaluation of Populations: A Simulation Study*. *Euphytica* **215**, 126. <https://doi.org/10.1007/s10681-019-2447-9>.
- van Frank G, Rivière P, Pin S, Baltassat R, Berthelot JF, Caizergues F, Dalmaso C, Gascuel JS, Hyacinthe A, Mercier F, Montaz H, Ronot B, Goldringer I (2020). *Genetic Diversity and Stability of Performance of Wheat Population Varieties Developed by Participatory Breeding*. *Sustainability* **12**, 384. <https://doi.org/10.3390/su12010384>.
- Vehtari A, Gelman A, Gabry J (2017). *Practical Bayesian Model Evaluation Using Leave-One-out Cross-Validation and WAIC*. *Statistics and Computing* **27**, 1413–1432. <https://doi.org/10.1007/s11222-016-9709-3>.
- Virk DS, Chakraborty M, Ghosh J, Prasad SC, Witcombe JR (2005). *Increasing the Client Orientation of Maize Breeding Using Farmer Participation in Eastern India*. *Experimental Agriculture* **41**, 413–426. <https://doi.org/10.1017/S001447970500270X>.
- Walter É, Pronzato L (1997). *Identification of parametric models from experimental data*. Communications and control engineering. London: Springer Masson.
- Welham SJ, Gogel BJ, Smith AB, Thompson R, Cullis BR (2010). *A comparison of analysis methods for late-stage variety evaluation trials*. *Australian & New Zealand Journal of Statistics* **52**, 125–149. <https://doi.org/10.1111/j.1467-842x.2010.00570.x>.
- Wolfe MS, Baresel JP, Desclaux D, Goldringer I, Hoad S, Kovacs G, Löschenberger F, Miedaner T, Østergård H, Lammerts van Bueren ET (2008). *Developments in Breeding Cereals for Organic Agriculture*. *Euphytica* **163**, 323–346. <https://doi.org/10.1007/s10681-008-9690-9>.

- Woyann LG, Benin G, Storck L, Trevizan DM, Meneguzzi C, Marchioro VS, Tonatto M, Madureira A (2017). *Estimation of Missing Values Affects Important Aspects of GGE Biplot Analysis*. *Crop Science* **57**, 40–52. <https://doi.org/10.2135/cropsci2016.02.0100>.
- Wricke G (1962). *Über Eine Methode Zur Erfassung Der Ökologischen Streubreite in Feldversuchen*. *Z. Pflanzenzuchtg* **47**, 92–96. <https://doi.org/10.1007/bf01103423>.
- Yan W (2013). *Biplot Analysis of Incomplete Two-way Data*. *Crop Science* **53**, 48–57. <https://doi.org/10.2135/cropsci2012.05.0301>.
- Yates F, Cochran WG (1938). *The Analysis of Groups of Experiments*. *The Journal of Agricultural Science* **28**, 556–580. <https://doi.org/10.1017/S0021859600050978>.