



HAL
open science

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

Michel Turbet Delof, Pierre Rivière, Julie C Dawson, Arnaud Gauffreteau, Isabelle Goldringer, Gaëlle van Frank, Olivier David

► To cite this version:

Michel Turbet Delof, Pierre Rivière, Julie C Dawson, Arnaud Gauffreteau, Isabelle Goldringer, et al.. Bayesian joint-regression analysis of unbalanced series of on-farm trials. 2024. hal-04380787v2

HAL Id: hal-04380787

<https://hal.science/hal-04380787v2>

Preprint submitted on 30 Sep 2024 (v2), last revised 12 Nov 2024 (v3)

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.



Distributed under a Creative Commons Attribution - NonCommercial 4.0 International License

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

Michel Turbet Delof¹, Pierre Rivière², Julie C. Dawson³, Arnaud Gauffreteau⁴, Isabelle Goldringer¹, Gaëlle van Frank¹, and Olivier David⁵

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 \times environment ($G \times E$) 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 $G \times E$ combinations from 189 on-farm trials. This dataset was severely unbalanced with more than 90% of $G \times E$ 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 as 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.

Keywords: decentralized participatory plant breeding; bread wheat; $G \times E$ interaction; hierarchical model; Finlay-Wilkinson model; Student's t distribution; varietal stability

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

Correspondence: michel.turbet_delof@cirad.fr

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 according to their own logic. 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

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

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

69 As very few varieties were common to all the trials and many varieties were tested in a lim-
70 ited number of trials, the resulting series of trials was very unbalanced, so that the estimation
71 of variety average performances and stabilities was difficult. Joint regression is a robust method
72 for estimating genetic main effects and stability with incomplete datasets (Finlay and Wilkin-
73 son, 1963; Pereira et al., 2007; Yates and Cochran, 1938). It is based on the Finlay-Wilkinson
74 (FW) model, which is parsimonious since the interaction effect between a genotype and an en-
75 vironment is modelled as the product of a genotype stability parameter, called sensitivity, and
76 the environment main effect. Various Finlay-Wilkinson models have been used in a frequen-
77 tist framework, in which environment effects were either fixed or random (Nabugoomu et al.,
78 1999; Ng and Williams, 2001; Patterson and Silvey, 1980). In the latter case, environment ef-
79 fects were assumed to come from a common distribution, thereby leading to shrunk estimates.
80 FW models in which genetic main effects, environment main effects and genetic sensitivities
81 (FW coefficient of regression) are all random effects have recently been developed. These have
82 been implemented in a Bayesian framework and when they include random effects, these are
83 called hierarchical models (Carlin and Louis, 2008; Robert, 2007). Thus far, these models have
84 been used to analyze slightly unbalanced trials (Lian and de los Campos, 2016). Hierarchical joint
85 regression has also been used to analyze very unbalanced simulated data (van Frank et al., 2019).
86 This simulation study has shown that genotypes should be tested in sufficiently many trials in
87 order to estimate their main effects and sensitivities reliably. However, this method had not been
88 used to analyze real and very unbalanced trials. Thus, it was not clear if it could cope with the
89 actual levels of unbalanced data seen in the French PPB on-farm trials and what insight it could
90 give into the behavior of genotypes across environments.

91 Extreme data is an important issue in data analysis. In multi-environment trials (MET), they
92 may come from either (1) errors between scoring and data formatting (measurement error, wrong
93 labelling, etc.), or (2) environmental heterogeneity in the trial (weed infestation, soil fertility, etc.),
94 or (3) the heterogeneity of the responses of the varieties tested between trials ($G \times E$ interaction).
95 In our PPB program, as cultivation environments are less controlled, extreme observations (types
96 2 and 3) could be more frequent than expected. This could reduce the precision of estimates
97 based on the normal distribution. Extreme observations could be removed from the dataset to
98 solve this problem, but it is difficult to decide which observations to remove. If too many ex-
99 treme observations are removed, then the variability of the data may be underestimated and
100 the precision of the statistical analysis overestimated. Alternatively, statistical methods that are
101 robust to extreme observations may be used (Hampel et al., 2011; Huber and Ronchetti, 1981).
102 Various robust methods have been developed in a frequentist or a Bayesian framework, in par-
103 ticular methods consisting in replacing the normal distribution by a Student's t distribution in
104 statistical models. This distribution is more robust to extreme observations than the normal dis-
105 tribution, because it has heavier tails (Carlin and Polson, 1991; Choy and Chan, 2008; Lange
106 et al., 1989; Rosa et al., 2003). It has been used to handle the extreme observations of a single

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

109 This study was aimed at developing statistical methods for analyzing series of on-farm trials,
 110 and at improving the assessment of varieties of the wheat PPB program by using the information
 111 at the level of the network. As our dataset was very unbalanced and could include extreme
 112 observations, we compared several Finlay-Wilkinson models, in particular hierarchical models
 113 and models based on the t distribution. These models were developed in a Bayesian framework,
 114 since this framework is rigorous and since it facilitates the implementation of complex models
 115 (Carlin and Louis, 2008; Robert, 2007). Finally, the best Finlay-Wilkinson model we obtained was
 116 used to analyze our data and characterize the behaviour of our varieties across environments.

117

2. Materials and methods

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
α	Germplasm main effect
θ	Environment main effect
η	Germplasm sensitivity (FW coefficient)
S^2	Germplasm static stability
W	Germplasm ecovalence (a dynamic stability)
LOO	Leave one out
elpd_{loo}	LOO expected logarithmic predictive density

Table 1 – Main notations.

118 In our study, a population variety is defined as a set of individuals which may be different
 119 but which are derived from the using certain agronomic practices, and a germplasm as any bio-
 120 logical entity whose individuals are derived from the same breeding process, including varieties
 121 registered in the official catalog, landraces, historic varieties, mixtures or populations stemming
 122 from crosses. An environment is the combination of a farm and a year.

2.1. Statistical methods

124 2.1.1. *Models.* We consider methods for analyzing series of on-farm trials in two steps (Patter-
 125 son, 1997; Patterson and Silvey, 1980). First, germplasm means are estimated using within-trial
 126 analyses, taking into account any block effects (spatial effects). Then, these estimates are ana-
 127 lyzed using a between-trial analysis. In the between-trial analysis, the phenotypic value $Y_{ij} \in \mathbb{R}$
 128 for a given trait Y , germplasm i and environment j is assumed to be equal to

$$129 \quad Y_{ij} = \mu_{ij} + \varepsilon_{ij},$$

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

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

$$134 \quad \mu_{ij} = \alpha_i + \theta_j,$$

135 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 .
 136 Models FWHs, FWs and FWHn model $G \times E$ interactions using the Finlay-Wilkinson regression,

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

Table 2 – The five models compared.

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

$$139 \quad \mu_{ij} = \alpha_i + \theta_j + \eta_i\theta_j,$$

140 where $\eta_i \in \mathbb{R}$ is the sensitivity of germplasm i to environments (linear regression coefficient,
141 Perkins and Jinks, 1968). As the average sensitivity is equal to 0, a germplasm with $\eta_i > 0$ is
142 more sensitive and germplasm with $\eta_i < 0$ is less sensitive to environments than a germplasm
143 with the average sensitivity. In these models, a part of the interaction between germplasm i
144 and environment j is modelled as a multiplicative term $\eta_i\theta_j$. The Finlay-Wilkinson coefficient is
145 considered as both a static and a dynamic indicator of stability (Becker and Leon, 1988; Lin et al.,
146 1986). In this model, statically stable genotypes have a coefficient close to -1. Dynamically stable
147 genotypes have a coefficient close to zero, but having a coefficient close to zero is not sufficient
148 to determine dynamic stability, this also depends on the amount of $G \times E$ variation that remains
149 unexplained by the model.

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

$$155 \quad \varepsilon_{ij} \sim \mathcal{N}\left(0, \sigma_\varepsilon^2\right),$$

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

$$160 \quad \varepsilon_{ij} \sim t\left(0, \sigma_\varepsilon^2, \nu\right),$$

161 where $t(0, \sigma_\varepsilon^2, \nu)$ is the Student's t distribution with dispersion parameter $\sigma_\varepsilon^2 > 0$ and $\nu > 2$
162 degrees of freedom. We assume that $\nu > 2$ to ensure that the expectation and the variance
163 of ε_{ij} are defined and finite. In models FWHs, FWs and ADHs, the variance of ε_{ij} is equal to
164 $\nu\sigma_\varepsilon^2/(\nu - 2)$. The normal distribution can be considered as a t distribution with ν tending to $+\infty$.
165 For additive models, the between-trial residuals combine the $G \times E$ effects and within-trial errors,
166 i.e. experimental errors and environmental heterogeneity in each trial, while for FW models, they
167 combine the part of $G \times E$ effects not explained by the multiplicative term $\eta_i\theta_j$ and within-trial
168 errors. Student residuals better handle data heterogeneity than normal residuals, since they can
169 be written as (Simar, 2002)

$$170 \quad \varepsilon_{ij} \sim \mathcal{N}\left(0, \sigma_{ij}^2\right), \quad \sigma_{ij}^{-2} \sim \Gamma(\nu/2, \nu\sigma_\varepsilon^2/2),$$

171 where $\Gamma(\nu/2, \nu\sigma_\varepsilon^2/2)$ is the gamma distribution with shape parameter $\nu/2$ and rate parameter
172 $\nu\sigma_\varepsilon^2/2$.

173 **2.1.2. Prior distribution.** The statistical methods are implemented in a Bayesian framework, so
174 that a joint prior distribution is placed on model parameters. Weakly informative prior distribu-
175 tions are placed on σ_ε and ν (Cao et al., 2022; Gelman, 2006; Juárez and Steel, 2010):

$$176 \quad \sigma_\varepsilon \sim \mathcal{N}^+(0, \lambda_\varepsilon^2), \quad \nu = 2 + \gamma, \quad \gamma \sim \Gamma(2, 0.1),$$

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

179 Since series of on-farm trials are often unbalanced and often involve many germplasm and
180 environments, α_i , θ_j and when present η_i are assumed to follow hierarchical distributions in all
181 the models except model FWs:

$$182 \quad \alpha_i \sim \mathcal{N}(\mu_Y, \sigma_\alpha^2), \quad \eta_i \sim \mathcal{N}(0, \sigma_\eta^2), \quad \theta_j \sim \mathcal{N}(0, \sigma_\theta^2),$$

183 where μ_Y , σ_α , σ_η and σ_θ are unknown parameters. Then, weakly informative prior distributions
184 are placed on the hyperparameters μ_Y , σ_α , σ_η and σ_θ :

$$185 \quad \mu_Y \sim \mathcal{N}(\lambda_\mu, \lambda_\varepsilon^2), \quad \sigma_\alpha \sim \mathcal{N}^+(0, \lambda_\varepsilon^2), \quad \sigma_\theta \sim \mathcal{N}^+(0, \lambda_\varepsilon^2), \quad \sigma_\eta \sim \mathcal{N}^+(0, 0.75^2),$$

186 where λ_μ is a known prior value of the trait mean. Germplasm main effects, environment main
187 effects, germplasm sensitivities and residuals are assumed to be independent given the hyper-
188 parameters, σ_ε and ν . In model FWs, the hierarchical distributions of α_i , η_i and θ_j are replaced
189 by weakly informative prior distributions:

$$190 \quad \alpha_i \sim \mathcal{N}(\mu_Y, \lambda_\varepsilon^2), \quad \eta_i \sim \mathcal{N}(0, 0.75^2), \quad \theta_j \sim \mathcal{N}(0, \lambda_\varepsilon^2).$$

191 The values chosen for λ_ε and λ_μ are in Appendix A.1.

192

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

196 **2.1.3. Posterior distribution.** Bayesian inference is based on the posterior distribution of model
197 parameters. This distribution is estimated using Markov chain and Monte Carlo (MCMC) meth-
198 ods. These methods simulate the values of model parameters according to a Markov chain
199 that converges to the posterior distribution of these parameters (Robert, 2007). They are im-
200 plemented using R (R Core Team, 2014) and the package `rstan` (Stan Development Team,
201 2016), that performs Hamiltonian Monte Carlo (HMC) sampling. This method aims at reducing
202 the correlation between successive sampled values by using a proposal distribution based on
203 Hamiltonian dynamics (Neal, 2011).

204 **2.1.4. Model comparison.** The predictive ability of models is compared using leave-one-out cross-
205 validation, which seems more appropriate than Bayes factors for selecting models that approx-
206 imate the process generating the data (Lartillot, 2023). We estimate the expected logarithmic
207 predictive density using the R package `LDD` (Vehtari et al., 2017). This criterion is equal to

$$208 \quad \text{elpd}_{\text{loo}} = \sum_{(i,j) \in \mathcal{C}} \ln(p(Y_{ij} | Y_{-ij})),$$

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

214 For main effects and sensitivities, we estimate the average standard deviation of estimates,
215 which allows us to estimate the precision of the analysis. To be able to compare the precision
216 between traits, for α and θ we estimate the average coefficient of variation by dividing this
217 standard deviation by the general average μ_Y .

218 **2.1.5. Variance decomposition.** In order to assess the importance of model terms, the variance of
219 an observation is decomposed for the main model FWs. Since α_i , θ_j , η_i and ε_{ij} are conditionally
220 independent, the terms θ_j and $\eta_i\theta_j$ are not correlated, and the variance of an observation given
221 the hyperparameters, σ_ε^2 and ν is equal to

$$222 \quad \text{Var}(Y_{ij}) = \text{Var}(\alpha_i + \theta_j + \eta_i\theta_j + \varepsilon_{ij}) = \sigma_\alpha^2 + \sigma_\theta^2 + \sigma_\eta^2\sigma_\theta^2 + \text{Var}(\varepsilon_{ij}).$$

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

$$226 \quad \pi(\alpha) = \frac{\sigma_\alpha^2}{\text{Var}(Y_{ij})}, \quad \pi(\theta) = \frac{\sigma_\theta^2}{\text{Var}(Y_{ij})}, \quad \pi(\eta\theta) = \frac{\sigma_\eta^2\sigma_\theta^2}{\text{Var}(Y_{ij})}.$$

227 $\pi(\alpha)$ is also called broad-sense heritability. The proportion of variance explained by the model
 228 (coefficient of determination) is equal to

$$229 \quad R^2 = \pi(\alpha) + \pi(\theta) + \pi(\eta\theta) = \frac{\sigma_\alpha^2 + \sigma_\theta^2 + \sigma_\eta^2\sigma_\theta^2}{\text{Var}(Y_{ij})}.$$

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

$$232 \quad \rho = \frac{\text{Var}(\eta_i\theta_j)}{\text{Var}(\eta_i\theta_j + \varepsilon_{ij})} = \frac{\sigma_\eta^2\sigma_\theta^2}{\sigma_\eta^2\sigma_\theta^2 + \text{Var}(\varepsilon_{ij})}.$$

233 **2.1.6. Characterization of germplasm.** The main effect and sensitivity of each germplasm are esti-
 234 mated using model FWHs. In addition, two stability indicators are estimated for each germplasm,
 235 the static stability S_i^2 (Becker and Leon, 1988) and the ecovalence W_i (Wricke, 1962) which is an
 236 indicator of dynamic stability. Due to data imbalance, the empirical estimates of these indicators
 237 are biased. Thus, we define stability indicators by means of theoretical variances using model
 238 FWHs (Cotes et al., 2006; Piepho, 1999). Using the independence assumptions of the model,
 239 we obtain for germplasm i ,

$$240 \quad W_i = \text{Var}(\eta_i\theta_j + \varepsilon_{ij}) = \eta_i^2\sigma_\theta^2 + \text{Var}(\varepsilon_{ij}),$$

$$241 \quad S_i^2 = \text{Var}(\theta_j + \eta_i\theta_j + \varepsilon_{ij}) = (1 + \eta_i)^2\sigma_\theta^2 + \text{Var}(\varepsilon_{ij}) = (1 + 2\eta_i)\sigma_\theta^2 + W_i.$$

242 The larger these indicators, the less stable the germplasm. Becker (1981) applied the same de-
 243 composition with the empirical variances.

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

250 2.2. Wheat PPB program

251 **2.2.1. Germplasm.** We studied 206 germplasm covering different "germplasm types": 98 "cross"
 252 germplasm resulting from crosses made either on the farm or at the research station (Rivière et
 253 al., 2015b), 50 "landraces", i.e. population varieties grown before 1884 (date of creation of Dattel,
 254 the first wheat variety from a controlled cross), 30 "historic varieties", developed by professional
 255 breeding before 1950, 17 "mixtures", which were generally complex, with numerous genotypes
 256 from potentially all the other germplasm types. In addition, 11 "registered varieties" after 1950
 257 and widely used in organic farming were included: Maitre Pierre (1954), Poncheau (1956), Renan
 258 (1990), Ataro (2004), Pollux (2004), Rubisco (2012), Hendrix (2012), Kampmann selected from
 259 Renan, and Hermes (1982), Alauda (2004) and Goldritter (2013), all three selected from Probus
 260 (1957).

261 **2.2.2. Experimental designs.** The data analyzed were collected between 2008 and 2019. The
 262 wheat PPB program followed numerous experimental designs due to the different constraints
 263 of farmers, collectives and researchers. The designs have been grouped into three classes (Tab.
 264 3). Some experimental designs (without blocks with repeated germplasm, and incomplete blocks
 265 with two blocks) were co-designed to be adapted to breeders' objectives, farmers' constraints
 266 and agricultural routines (Dawson et al., 2011). In these designs, the germplasm common to
 267 all farms (control germplasm) were collectively chosen by farmers and researchers, while each
 268 farmer individually chose the additional germplasm to be cultivated in his farm. At the beginning

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

Designs	Nb of blocks	Nb of repeated germplasm	Nb of gemplasm by environment	Nb of envi- ronments
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

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

276 **2.2.3. Data collected.** Four traits were studied, plant height (60% of the data was the average
 277 height of 25 individuals and 40% was the overall height of the microplot, mm), spike weight (mean
 278 of 25 individual measures, g), protein content of the grain (on the microplot, measured with NIRS
 279 technology at INRAE Clermont-Ferrand France, %) and thousand kernel weight (TKW, measured
 280 on the microplot, g). These four traits were among those collectively chosen by farmers and
 281 researchers to be measured during the PPB program (Tab. 4). Plant height was measured in the
 282 field, while the other traits were measured after harvest at the research station on samples of
 283 spikes sent by farmers. Outliers with respect to agronomic knowledge of the traits were excluded
 284 (for example, a plant taller than three meters).

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

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

297 These data were analyzed using the models of Tab. 2. As the dataset was very unbalanced, it
 298 was not clear if model parameters were identifiable. Thus, for each variable, the identifiability of
 299 germplasm main effects and environment main effects was studied for the additive model. We
 300 checked that the rank of the design matrix of the model was equal to $1 + (I - 1) + (J - 1)$, where
 301 I was the number of germplasm and J the number of environments (p. 50, Silvey, 1975). For the
 302 FW model, identifiability was more difficult to study because the model was nonlinear. Thus, we
 303 restricted ourselves to studying local identifiability near an estimate of model parameters (Chap.
 304 2, Walter and Pronzato, 1997). First, a linear approximation of the model was carried out using a
 305 Taylor expansion. Then, we checked that the rank of the design matrix of this linear model was
 306 equal to $1 + (I - 1) + (J - 1) + (I - 1)$.

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

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

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

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

314

3. Results

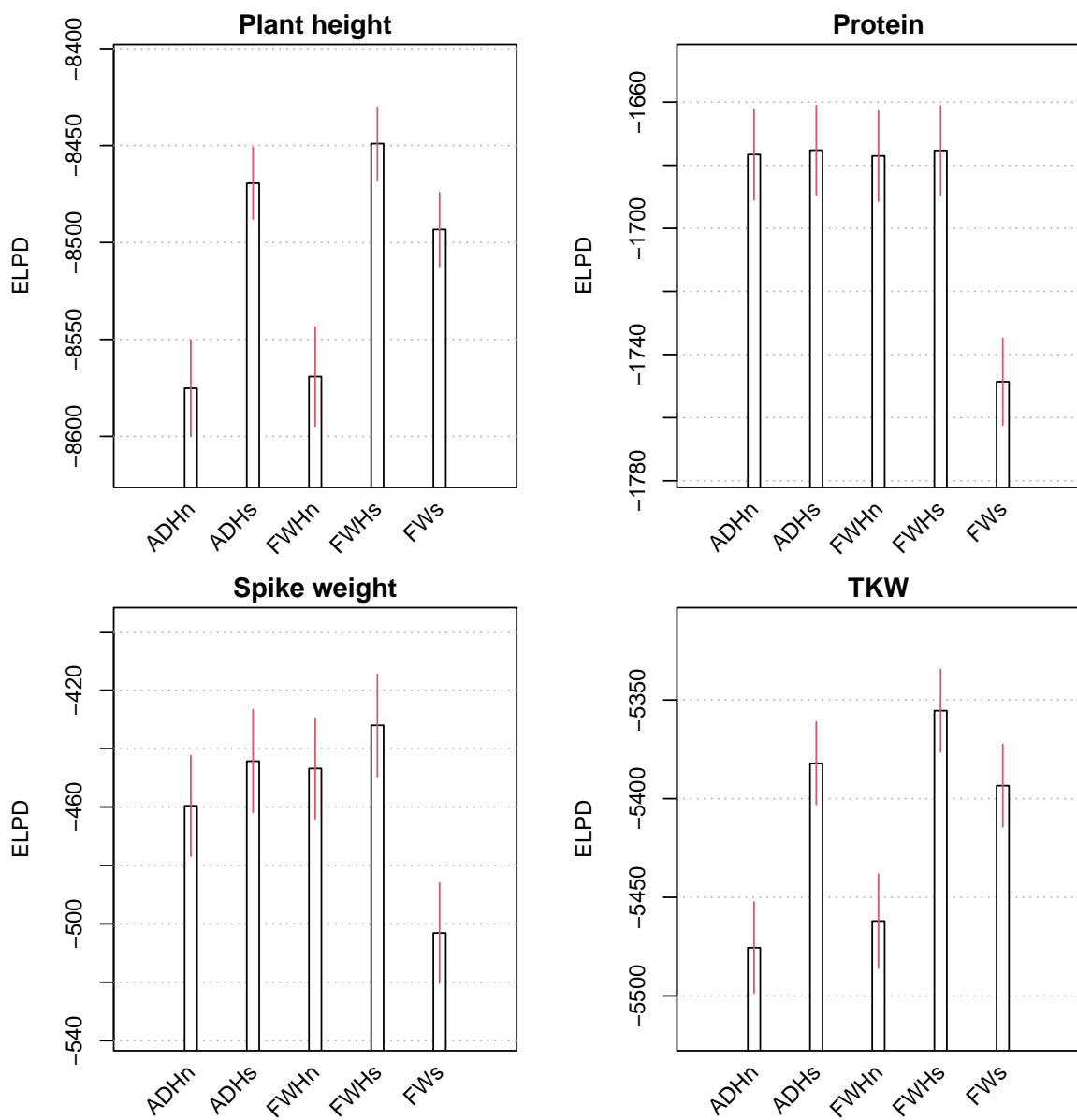


Figure 1 - Predictive capacity of models. elpd₁₀₀ and its associated standard error for the four studied traits.

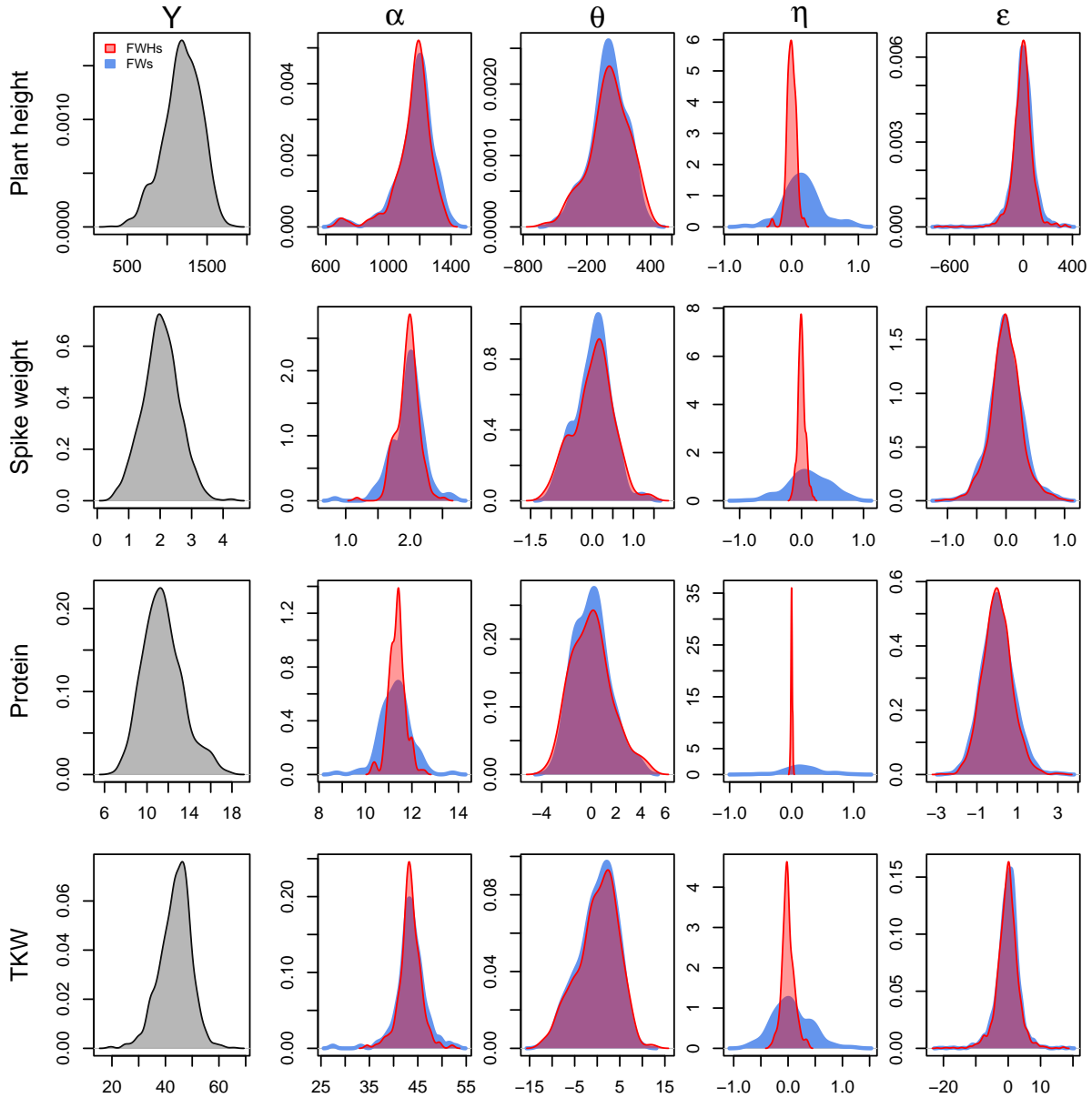


Figure 2 – The first column presents the distribution of the trait to be explained (in grey). The last four columns compare the hierarchical (red) and non hierarchical (blue) versions of the FW model with a Student law for the residuals, and show the smoothed histograms of main effects, FW coefficients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

315 3.1. Predictive capacity of models

316 According to the $elpd_{loo}$ criterion, the non-hierarchical FWs model was less predictive than
 317 the hierarchical FWs model for all the traits (Fig. 1). Using the latter model shrank the estimates
 318 of η and sometimes α (Fig. 2). With the non-hierarchical model (FWs), some estimates (α_i and
 319 η_j) seemed to be unreliable, in particular some germplasm means were extreme and some FW
 320 coefficients were larger than 1 or smaller than -1.

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

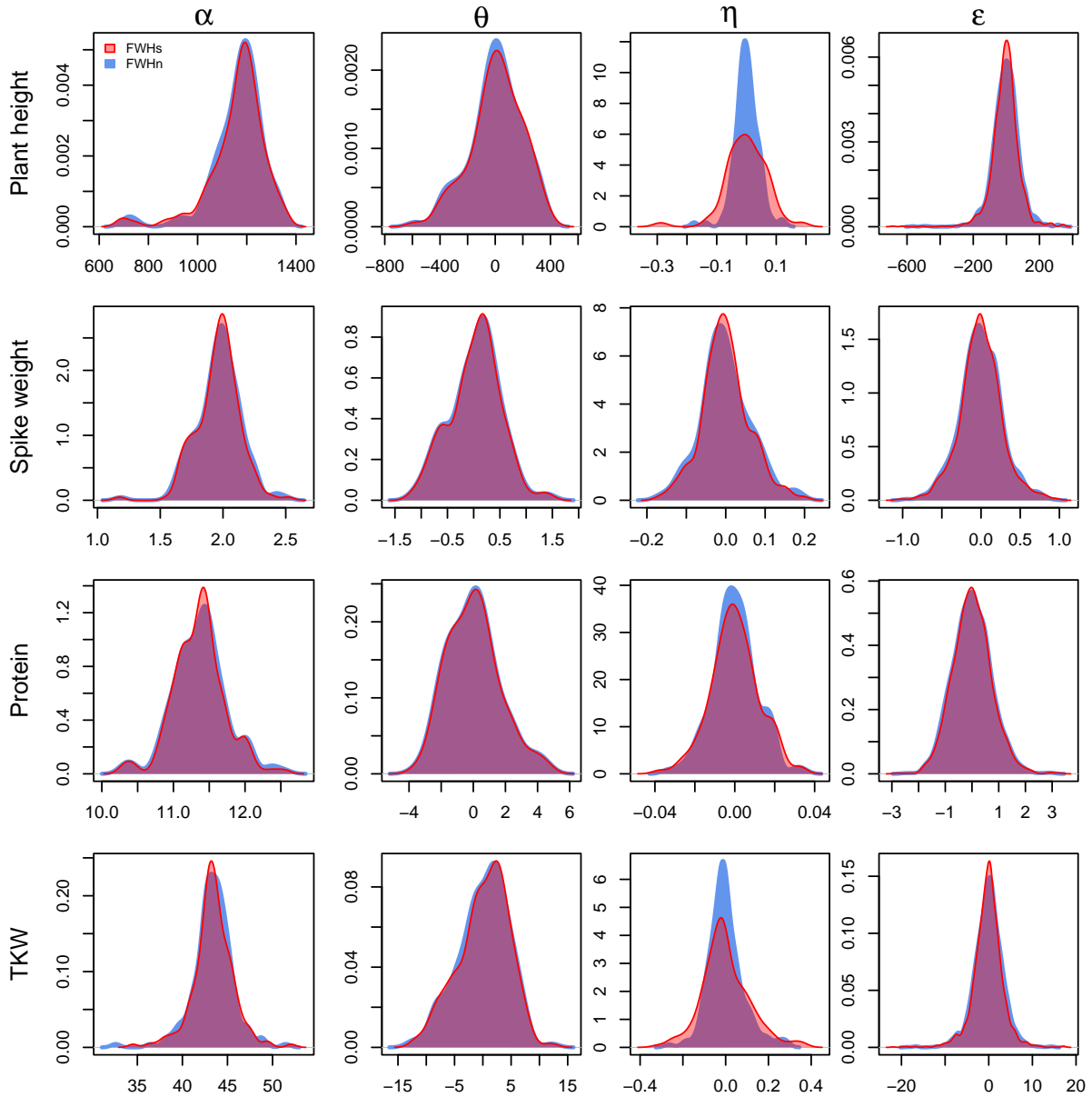


Figure 3 – Comparison of hierarchical FW models with different residual laws, the Student (red) and the normal (blue). These graphics show the smoothed histograms of main effects, FW coefficients and residuals.

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

331 The $elpd_{100}$ criterion was estimated using Pareto smoothed importance sampling (Vehtari et
 332 al., 2017). This method tends to be less precise for models that do not fit the data well. Thus,
 333 as expected, estimates of $elpd_{100}$ were more reliable for the two hierarchical models with a t
 334 likelihood (FWHs and ADHs) than for the other models, in particular model FWs (Supplementary
 335 Tab. B.1).

336 3.2. Precision of estimates and distribution of residuals

337 For the models with a t distribution, the estimate of the number of degrees of freedom (ν)
 338 varied between 3.4 and 27.6 (close to a normal distribution) (Tab. 5). Thus, the shape of the dis-
 339 tribution of residuals depended on the trait. This result confirmed that the number of extreme

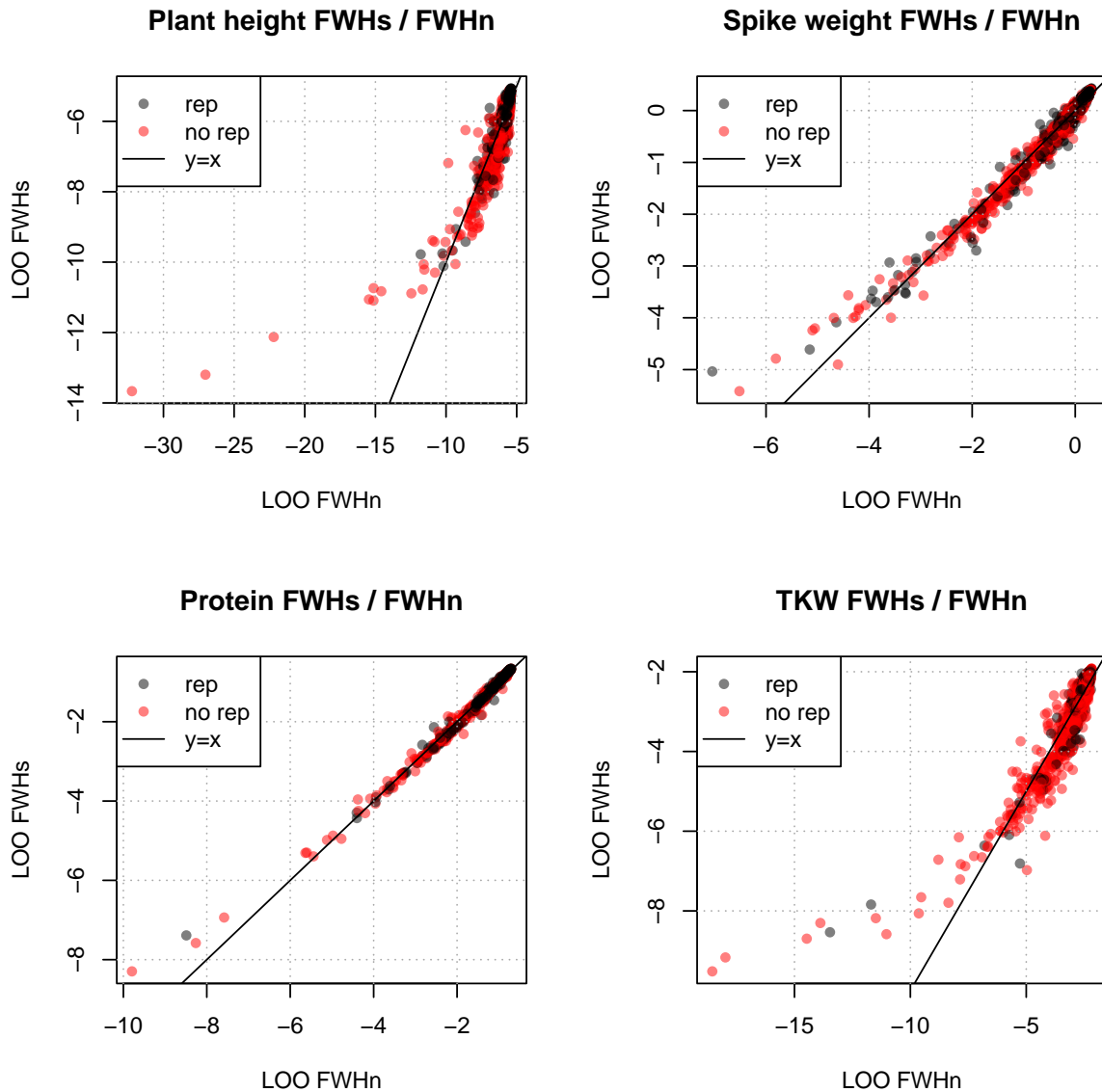


Figure 4 – Comparison of t and normal models (FWHs vs FWHn) in terms of the contributions of observations to the elpd_{loo} criterion. Black (resp. red) dots correspond to observations that were measured on germplasm that were replicated (resp. not replicated) within trials.

340 observations was not negligible in our data, and that models with a t distribution were more ap-
 341 propriate. In the latter case, the variation ranges of residuals were wider but with more residual
 342 values close to 0 for the t distribution than the normal distribution (Fig. 3). Models had similar es-
 343 timated precision, except for model FWs, which had less precise estimates. This result confirmed
 344 that a basic joint regression, i.e. non-hierarchical model, was not suited to our unbalanced data.
 345 Parameters α and θ were estimated more precisely (difference in coefficient of variation between
 346 0 and 1.9, Tab. 5) for t models (ADHs and FWHs) than for normal models (ADHn and FWHn).
 347 This result was consistent with Fig. 4, where extreme observations were better predicted by
 348 more likely under model FWHs than under FWHn, except for protein content.

349 3.3. Variance decomposition

350 The proportion of variance explained by each term of model FWHs depended on the trait
 351 (Tab. 6). For all four traits, the environment effect was highly explanatory. For height and TKW,

Trait	Model	ν	$cv(\alpha)$	$cv(\theta)$	$sd(\eta)$
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 5 – Number of degrees of freedom and precision of estimates.

ν : 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).

	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]
ρ	6.7	[2.9, 12.1]	4.9	[1.7, 9.3]	1.2	[0, 4.6]	6.9	[3.1, 12]

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 proportion of variance explained by α , θ and $\eta\theta$. ρ is the proportion of the variance of $G \times E$ and errors explained by $\eta\theta$.

352 a relatively large part of the total variance was explained by the germplasm effect (resp. 24%
 353 and 16.1%), whereas this part was much smaller for spike weight and protein content (10.9%
 354 and 5.7%). The proportion of variance explained by the sensitivity effect η was not significantly
 355 different from 0 for protein content and low for the three other traits. It explained 6.7%, 4.9%
 356 and 6.9% of the variance of $G \times E$ interactions and experimental errors (ρ parameter) for plant
 357 height, spike weight and TKW, respectively.

358 3.4. Characterization of germplasm

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

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 7 – Correlation between germplasm parameters.

*, **, *** : significant at $P = 0.05$, $P = 0.01$, $P = 0.001$ respectively.

α_i : germplasm effect, η_i : germplasm sensitivity (FW coefficient), S_i^2 : static stability, W_i : ecovalence.

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

371 4. Discussion

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

376 4.1. Handling the data from a highly unbalanced series of trials

377 As the farmers of the program chose the germplasm they assessed, the data obtained from
 378 the series of trials were very unbalanced, with more than 90% of the $G \times E$ combinations miss-
 379 ing. This made the estimation of germplasm main effects and sensitivities difficult. Although the
 380 Finlay-Wilkinson model was parsimonious, a basic joint regression with weakly-informative prior
 381 distributions (model FWs) was not able to cope with this level of disequilibrium. According to the
 382 $elpd_{100}$ criterion, model FWs was not the best model (Fig. 1). In addition, its estimates had poor
 383 precision and it led to extreme sensitivity estimates, with values close to 1 or -1 (Fig. 2).

384 In contrast, hierarchical joint regression appeared more suited to our data structure. Model FWs
 385 had the largest $elpd_{100}$ values for three traits out of four. Placing a hierarchical distribution on
 386 sensitivities constrained estimates and brought them closer to 0. This led to more satisfactory
 387 sensitivity estimates, since they were well below 1 in absolute value.

388 Three strategies have previously been used to manage incomplete $G \times E$ data: i) subset the
 389 total dataset to obtain an almost balanced subset for the analysis (Ceccarelli and Grando, 2007),
 390 ii) predict missing data with a more or less complex model and use these predictions in the
 391 analysis (Kumar et al., 2012; Woyann et al., 2017), and iii) use a model more robust to unbalanced
 392 data, provided it complies with model validation conditions (Assis et al., 2018; van Frank et al.,
 393 2019). We used the last strategy to maximise the amount of information from the data (less data
 394 excluded than in the first strategy) with a one-step process (unlike the second strategy).

395 Cotes et al. (2006) used a Bayesian approach to estimate FW coefficients in a MET study in
 396 order to take prior information on germplasm coming from other studies into account. A similar
 397 approach was used by Couto et al. (2015), Foucteau and Denis (2001), and Nascimento et al.
 398 (2020) and was found to greatly improve the results. Here, we used little prior information. But
 399 in the future, previous evaluation studies may provide stronger prior information on germplasm
 400 behaviour.

Trait	Registered	Historic	Landrace	Cross	Mixture	
Plant height	$\bar{\alpha}_k$	862 ^d	1136 ^c	1220 ^a	1175 ^b	1188 ^b
		[822, 901]	[1096, 1175]	[1181, 1258]	[1138, 1210]	[1147, 1228]
	$\bar{\eta}_k$	-0.11 ^b	-0.01 ^a	0 ^a	0.01 ^a	-0.01 ^a
		[-0.2, -0.03]	[-0.06, 0.05]	[-0.05, 0.05]	[-0.01, 0.04]	[-0.09, 0.06]
	\bar{S}_k^2	37688 ^b	44351 ^{ab}	44813 ^a	45620 ^a	43800 ^{ab}
	[28956, 48901]	[34878, 56497]	[35456, 57056]	[36512, 57472]	[34101, 56150]	
	\bar{W}_k	9067 ^a	7959 ^b	7988 ^b	7880 ^b	7772 ^b
		[7272, 11725]	[6568, 10184]	[6593, 10234]	[6541, 10056]	[6419, 10000]
Spike weight	$\bar{\alpha}_k$	2.02 ^b	1.98 ^{ab}	1.93 ^a	1.96 ^a	2.01 ^b
		[1.92, 2.12]	[1.89, 2.08]	[1.85, 2.02]	[1.87, 2.04]	[1.92, 2.11]
	$\bar{\eta}_k$	0.02 ^a	0.02 ^a	-0.01 ^a	0 ^a	-0.01 ^a
		[-0.06, 0.1]	[-0.03, 0.08]	[-0.05, 0.03]	[-0.02, 0.03]	[-0.07, 0.05]
	\bar{S}_k^2	0.34 ^a	0.34 ^a	0.32 ^a	0.33 ^a	0.32 ^a
	[0.27, 0.42]	[0.28, 0.41]	[0.27, 0.39]	[0.28, 0.39]	[0.26, 0.39]	
	\bar{W}_k	0.08 ^a	0.08 ^a	0.08 ^a	0.08 ^a	0.08 ^a
		[0.08, 0.09]	[0.08, 0.09]	[0.08, 0.09]	[0.08, 0.09]	[0.08, 0.09]
Protein	$\bar{\alpha}_k$	11.08 ^a	11.37 ^b	11.4 ^b	11.35 ^b	11.47 ^b
		[10.73, 11.42]	[11.06, 11.7]	[11.09, 11.72]	[11.04, 11.66]	[11.13, 11.8]
	$\bar{\eta}_k$	0 ^a	0 ^a	0 ^a	0 ^a	0 ^a
		[-0.03, 0.04]	[-0.02, 0.02]	[-0.02, 0.02]	[-0.01, 0.01]	[-0.03, 0.03]
	\bar{S}_k^2	3.4 ^a	3.4 ^a	3.39 ^a	3.4 ^a	3.4 ^a
	[2.72, 4.29]	[2.73, 4.27]	[2.72, 4.26]	[2.73, 4.27]	[2.72, 4.28]	
	\bar{W}_k	0.62 ^a	0.61 ^a	0.62 ^a	0.62 ^a	0.61 ^a
		[0.56, 0.68]	[0.56, 0.68]	[0.56, 0.68]	[0.56, 0.68]	[0.56, 0.68]
TKW	$\bar{\alpha}_k$	43.6 ^{ab}	43.8 ^a	43.1 ^b	43.4 ^{ab}	43.9 ^a
		[42.6, 44.6]	[42.9, 44.8]	[42.3, 43.9]	[42.6, 44.1]	[43, 44.8]
	$\bar{\eta}_k$	-0.03 ^{ab}	0.08 ^a	-0.03 ^b	0.01 ^{ab}	-0.05 ^b
		[-0.14, 0.09]	[0, 0.17]	[-0.09, 0.03]	[-0.03, 0.05]	[-0.14, 0.04]
	\bar{S}_k^2	33.4 ^{ab}	37.8 ^a	33.3 ^b	34.9 ^{ab}	32.2 ^b
	[26.9, 41.4]	[31.6, 45.7]	[28.2, 39.7]	[29.7, 41.2]	[26.6, 39]	
	\bar{W}_k	13.2 ^a	13.4 ^a	13.4 ^a	13.2 ^a	13.1 ^a
		[11.5, 15.6]	[11.6, 15.7]	[11.7, 15.7]	[11.6, 15.5]	[11.4, 15.4]

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.

401 4.2. Extreme observations

402 Extreme observations were more frequent in our dataset than expected under the normal
403 distribution for three traits out of four (Fig. 4). For these traits, using a t distribution increased
404 elpd_{100} values, and the estimate of the number of degrees of freedom of this distribution was
405 smaller than 10 (Tab. 5). In our application, observations were germplasm means resulting from
406 within-trial analyses rather than plot measurements. Extreme observations could occur for sev-
407 eral reasons, for example because of the heterogeneity of within-trial residual variances and
408 replications, because cultivation environments were less controlled, or because a non-negligible
409 part of $G \times E$ interactions was not captured by the multiplicative term of the FW model. The nor-
410 mal distribution was appropriate for protein content. It is difficult to explain why this trait had
411 fewer extreme observations. A possible explanation could be that the measurement of protein
412 content is more standardized than other trait measurements. For plant height, extreme values oc-
413 curred only for non-replicated micro-plots with a global measurement and never with data from
414 the average of 25 plants (Sect. 2.2.3), suggesting that the plot measurement is less accurate.
415 For TKW, the kernel count could be affected by broken kernels due to over-drying or incorrect
416 threshing settings leading to an overestimation of the number of kernels in the sample. Another

417 possible explanation is that protein content is less variable under different conditions than plant
418 height and spike weight (Kazakou et al., 2014).

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

422 The Student distribution is expected to take better account of extreme data and to yield
423 more robust estimates (Besag and Higdon, 1999; Lange et al., 1989; Rosa et al., 2003). Extreme
424 data are more likely to occur when varieties are not replicated within trials, which is frequent
425 in this dataset (Fig. 4). Rosa et al. (2003) found that a normal likelihood misestimated a main
426 effect compared to a t likelihood. This effect was estimated less precisely with a normal distri-
427 bution, which is consistent with our results for plant height, spike weight and TKW. A Student
428 distribution appears to be a good solution for dealing with extreme data, in particular in stability
429 analyses, where extreme observations are sometimes removed (this is justified when they are
430 extreme because of experimental errors, but not when they are due to natural variability). While
431 this distribution has recently been used to implement robust alternatives to BLUP (Gianola et
432 al., 2018) or to handle environmental heterogeneity in a single trial (Cao et al., 2022), to our
433 knowledge, it has not already been used in MET studies.

434 4.3. Computing time

435 Series of trials often include many genotypes and environments, leading to large data sets.
436 Thus, their analysis using mixed or hierarchical models is generally computationally demanding
437 (Smith et al., 2005). The computational load can be reduced by using approximate estimation
438 methods (Nabugoomu et al., 1999) or efficient algorithms, such as algorithms based on sparse
439 matrix operations (Gilmour et al., 1995; Thompson et al., 2003). Hierarchical joint regression
440 has already been implemented using Gibbs sampling or Jags (Lian and de los Campos, 2016;
441 van Frank et al., 2020). Our implementation based on Hamiltonian Monte Carlo and Stan was
442 more efficient since it required fewer iterations. It allowed us to analyze large datasets in about
443 6 minutes.

444 To reduce computing time, the analyses were carried out in two steps. This two-stage ap-
445 proach analyzed $G \times E$ means without taking account of their standard error, which can reduce
446 the efficiency of the analysis (Welham et al., 2010; Yates and Cochran, 1938). It would be inter-
447 esting to develop a one-stage method for analyzing plot measurements, in order to better take
448 account of the heterogeneity of the within-trial residual variances and replications (Rivière et al.,
449 2015a).

450 4.4. Variance decomposition

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

458 This variance decomposition is useful to identify the model terms which are the most im-
459 portant. In our application, the environment effects were the most important, explaining from
460 51% for TKW to 77% for protein content of the variance of observations (Tab. 6). This result
461 is consistent with the diversity of the cropping environments encountered (soil, climate, crop-
462 ping practices...) and with previous studies (Lian and de los Campos, 2016; Patterson and Silvey,
463 1980; Talbot, 1984).

464 4.5. Germplasm main effects and stabilities

465 Heritability was significant with plant height > TKW > spike weight > protein. Rivière et al.
466 (2015b) found (with data included in our study) a similar ranking in heritability: plant height >
467 TKW = protein > spike weight. Plant height is known to be quite heritable due to a relatively

468 simple genetic architecture with a few major genes, such as the well known Green Revolution
469 Rht1 and Rht2 genes (Peng et al., 1999). In our study, the presence of both recently registered
470 varieties and varieties dating from before the second World War, very likely led to varieties
471 containing different alleles for these loci and increased variability for height. The decrease in
472 plant height from landraces to historic varieties and registered varieties appears very clearly
473 (Tab. 8) as also found in several studies (Bektas et al., 2016; Cantarel et al., 2021).

474 FW coefficients explained a low proportion of the total variance (between 0.2% and 2.2%)
475 and a low proportion of the variance of $G \times E$ interactions and errors (between 1.2% and 6.9%,
476 Tab. 6). We can presume that the explanation of the interaction by the FW parameter is weaker
477 the greater the number of environments, for example 29% with less than 10 environments (12
478 studies), and 12% with more than 10 environments (11 studies, Brancourt-Hulmel et al., 1997).
479 Other classical models, such as AMMI (additive main effect and multiplicative interaction) or
480 GGE ($G + G \times E$) models, might explain a larger part of $G \times E$ interactions. Missing data estimation
481 methods allow these models to be used when the data are highly unbalanced, with up to from
482 40% unbalanced data for a MET with less than 20 environments to 60% unbalanced data for MET
483 with at least 40 environments (Woyann et al., 2017; Yan, 2013). However, these datasets are
484 more balanced than ours, and, as found by Rodrigues et al. (2011), FW is more robust than AMMI
485 when the data are highly unbalanced (75%). In our study, most germplasm occurred in a limited
486 number of environments, so that a parsimonious and very simple modelling of $G \times E$ interactions
487 had to be used. An alternative approach would be to better characterize the environments and
488 thus explain the environmental effects and part of the $G \times E$ interaction using environmental
489 variables (Piepho and Blancon, 2023).

490 Although sensitivities explained a rather low proportion of variance, FWs model had larger
491 $elpd_{loo}$ values than additive models for three traits out of four. In addition, for these traits, some
492 sensitivity estimates were not negligible, with values close to 0.2 or 0.3. Interaction effects then
493 represented 20% or 30% of environmental effects. Additive models were appropriate for the
494 protein content trait. It was found that the multiplicative term of the FW model was not signifi-
495 cant for protein content, both in a balanced network of 15 environments in Serbia (Hristov et al.,
496 2010) and in 12 environments in Swiss organic trials (Knapp et al., 2017). On the contrary, Mut et
497 al. (2010) found significant FW coefficients for a balanced network of 7 environments in Turkey.
498 These contrasting results could be explained by differences between numbers of environments
499 or between genetic diversities.

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

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

512 Dependence between stability and mean is widespread (Reckling et al., 2021), but in our case,
513 the correlation was low, which simplified interpretation of the stability analysis. Several studies
514 for different traits and with balanced MET found a very strong correlation between FW coeffi-
515 cient and the static stability (Becker, 1981; Fasahat et al., 2015; Reckling et al., 2021). However,
516 in our case, this relationship was even stronger (Tab. 7), probably because of the assumption
517 that the variance of residuals did not depend on the genotype. As in many other studies, the
518 residual variance was assumed to be independent of germplasm throughout our study. Allow-
519 ing the residual variance to depend on the genotype could improve the estimates of stability
520 indicators (Cotes et al., 2006; Couto et al., 2015). In particular, the dynamic indicator would be
521 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).

Supplementary information

Appendix A. Models

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

	λ_{μ}	λ_{ε}	μ_{emp}	σ_{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

Table A.1 – Known values of the parameters of the prior distribution (λ_{μ} , λ_{ε}), empirical mean (μ_{emp}) and standard deviation (σ_{emp}) of traits.

Appendix B. Model comparison

561

562 Tab. B.1 provides supplementary information on the estimation of the elpd_{loo} criterion. Fig. B.1
 563 provides supplementary information on the comparison of models FWHs and ADHs. Fig. B.2
 564 provides supplementary information on the comparison of models ADHs and ADHn.

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

Table B.1 – Estimates of tail shape parameters (k) used to estimate elpd_{loo} . The contribution of each observation to elpd_{loo} , 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.

565

Acknowledgements

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

571

Fundings

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

575

Conflict of interest disclosure

576 The authors declare that they comply with the PCI rule of having no financial conflicts of
 577 interest in relation to the content of the article. [IF APPROPRIATE: The authors declare the
 578 following non-financial conflict of interest: XXX (if some of the authors are recommenders of a
 579 PCI, indicate it here)].

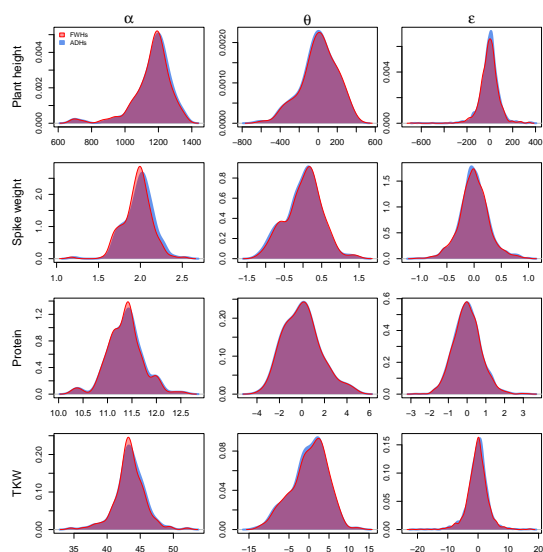


Figure B.1 – Comparison of models FWHs and ADHs for the distribution of germplasm main effects (α), environment main effects (θ) and residuals (ε) for each trait. Red: model FWHs; Blue: model ADHs.

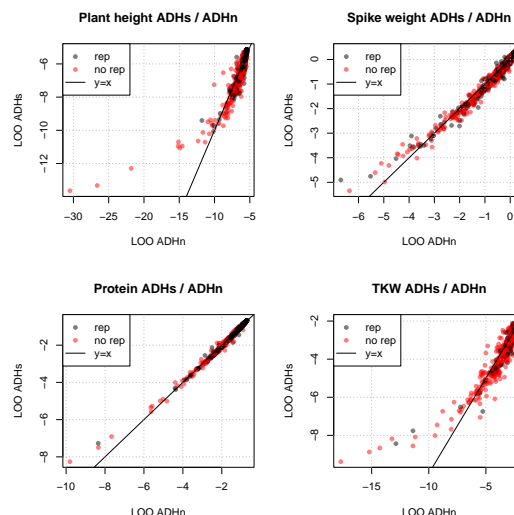


Figure B.2 – Comparison of models ADHs and ADHn in terms of the contributions of observations to the $elpd_{loo}$ criterion. Black (resp. red) dots correspond to observations that were measured on germplasm that were repeated (resp. not repeated) within trials.

Data, script, code, and supplementary information availability

580

581 Data and script for models are available online: DOI 10.57745/SUTZ9U <https://doi.org/10.57745/SUTZ9U>.

583

References

- 584 Annicchiarico P (2002). *Genotype X Environment Interactions: Challenges and Opportunities for Plant Breeding and Cultivar Recommendations*. Food & Agriculture Org.
- 585 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>.
- 586 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.
- 587 Assis TOG, Dias CTdS, Rodrigues PC (2018). *A Weighted AMMI Algorithm for Nonreplicated Data*. *Pesquisa Agropecuária Brasileira* **53**, 557–565.
- 588 Becker HC (1981). *Correlations among Some Statistical Measures of Phenotypic Stability*. *Euphytica* **30**, 835–840.
- 589 Becker HC, Leon J (1988). *Stability Analysis in Plant Breeding*. *Plant breeding* **101**, 1–23.
- 590 Bektas H, Hohn CE, Waines JG (2016). *Root and Shoot Traits of Bread Wheat (Triticum Aestivum L.) Landraces and Cultivars*. *Euphytica* **212**, 297–311.
- 591 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>.
- 592 Beza E, Steinke J, Van Etten J, Reidsma P, Fadda C, Mittra 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.

604

- 605 Brac de la Perrière RA, De Kochko P, Neubauer C, Storup B (2011). *Visions Paysannes de La*
606 *Recherche Dans Le Contexte de La Sélection Participative*.
- 607 Brancourt-Hulmel M, Biarnès-Dumoulin V, Denis JB (1997). *Points de repère dans l'analyse de la*
608 *stabilité et de l'interaction génotype-milieu en amélioration des plantes*. *Agronomie* **17**, 219–246.
609 <https://doi.org/10.1051/agro:19970403>.
- 610 Cantarel AA, Allard V, Andrieu B, Barot S, Enjalbert J, Gervais J, Goldringer I, Pommier T, Saint-
611 Jean S, Le Roux X (2021). *Plant Functional Trait Variability and Trait Syndromes among Wheat*
612 *Varieties: The Footprint of Artificial Selection*. *Journal of Experimental Botany* **72**, 1166–1180.
- 613 Cao Z, Stefanova K, Gibberd M, Rakshit S (2022). *Bayesian Inference of Spatially Correlated Ran-*
614 *dom Parameters for On-Farm Experiment*. *Field Crops Research* **281**, 108477.
- 615 Carlin BP, Polson NG (1991). *Inference for Nonconjugate Bayesian Models Using the Gibbs Sampler*.
616 *Canadian Journal of statistics* **19**, 399–405.
- 617 Carlin BP, Louis TA (2008). *Bayesian Methods for Data Analysis*. CRC Press.
- 618 Ceccarelli S, Grando S, Bailey E, Amri A, El-Felah M, Nassif F, Rezgui S, Yahyaoui A (2001). *Farmer*
619 *Participation in Barley Breeding in Syria, Morocco and Tunisia*. *Euphytica* **122**, 521–536. <https://doi.org/10.1023/A:1017570702689>.
- 620 Ceccarelli S, Grando S (2007). *Decentralized-Participatory Plant Breeding: An Example of Demand*
621 *Driven Research*. *Euphytica* **155**, 349–360. <https://doi.org/10.1007/s10681-006-9336-8>.
- 622 Ceccarelli S, Grando S (2020). *Participatory Plant Breeding: Who Did It, Who Does It and Where?*
623 *Experimental Agriculture* **56**, 1–11.
- 624 Choy BST, Chan JSK (2008). *Scale Mixtures Distributions in Statistical Modelling*. *Australian & New*
625 *Zealand Journal of Statistics* **50**, 135–146.
- 626 Cotes JM, Crossa J, Sanches A, Cornelius PL (2006). *A Bayesian Approach for Assessing the Stability*
627 *of Genotypes*. *Crop Science* **46**, 2654–2665. [https://doi.org/10.2135/cropsci2006.04.](https://doi.org/10.2135/cropsci2006.04.0227)
628 [0227](https://doi.org/10.2135/cropsci2006.04.0227).
- 629 Couto MF, Nascimento M, do Amaral Jr AT, e Silva FF, Viana AP, Vivas M (2015). *Eberhart and*
630 *Russel's Bayesian Method in the Selection of Popcorn Cultivars*. *Crop Science* **55**, 571–577.
- 631 David O, van Frank G, Goldringer I, Rivière P, Turbet Delof M (2020). *Bayesian Inference of Natural*
632 *Selection from Spatiotemporal Phenotypic Data*. *Theoretical Population Biology* **131**, 100–109.
633 <https://doi.org/10.1016/j.tpb.2019.11.007>.
- 634 Dawson J, Murphy KM, Jones SS (2008). *Decentralized Selection and Participatory Approaches in*
635 *Plant Breeding for Low-Input Systems*. *Euphytica* **160**, 143–154. [https://doi.org/10.1007/](https://doi.org/10.1007/s10681-007-9533-0)
636 [s10681-007-9533-0](https://doi.org/10.1007/s10681-007-9533-0).
- 637 Dawson JC, Rivière P, Berthelot JF, Mercier F, Kochko P, Galic N, Pin S, Serpolay E, Thomas M,
638 Giuliano S, Goldringer I (2011). *Collaborative Plant Breeding for Organic Agricultural Systems in*
639 *Developed Countries*. *Sustainability* **3**, 1206–1223. <https://doi.org/10.3390/su3081206>.
- 640 Desclaux D, Nolot JM, Chiffolleau Y, Gozé E, Leclerc C (2008). *Changes in the Concept of Genotype*
641 *× Environment Interactions to Fit Agriculture Diversification and Decentralized Participatory Plant*
642 *Breeding: Pluridisciplinary Point of View*. *Euphytica* **163**, 533–546. [https://doi.org/10.](https://doi.org/10.1007/s10681-008-9717-2)
643 [1007/s10681-008-9717-2](https://doi.org/10.1007/s10681-008-9717-2).
- 644 Döring TF, Annicchiarico P, Clarke S, Haigh Z, Jones HE, Pearce H, Snape J, Zhan J, Wolfe MS
645 (2015). *Comparative Analysis of Performance and Stability among Composite Cross Populations,*
646 *Variety Mixtures and Pure Lines of Winter Wheat in Organic and Conventional Cropping Systems*.
647 *Field Crops Research* **183**, 235–245.
- 648 Falconer DS (1960). *Introduction to Quantitative Genetics*. Edinburgh/London: Olivier & Boyd.
- 649 Fasahat P, Rajabi A, Mahmoudi SB, Noghabi MA, Rad JM (2015). *An Overview on the Use of Sta-*
650 *bility Parameters in Plant Breeding*. *Biometrics & Biostatistics International Journal* **2**, 00043.
- 651 Finlay KW, Wilkinson GN (1963). *The Analysis of Adaptation in a Plant-Breeding Programme*. *Aus-*
652 *tralian Journal of Agricultural Research* **14**, 742–754. <https://doi.org/10.1071/AR9630742>.
- 653 Fouceteau V, Denis JB (2001). *Statistical Analysis of Successive Series of Experiments in Plant Breed-*
654 *ing: A Bayesian Approach*. *Quantitative genetics and breeding methods: the way ahead. Proceed-*
655 *ings of the Eleventh Meeting of the EUCARPIA Section Biometrics in Plant Breeding, Paris, France,*
656 *30/31 August - 1 September, 2000*, 49–56.
- 657

- 658 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>.
- 659
- 660
- 661 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>.
- 662
- 663
- 664 Gelman A, Goodrich B, Gabry J, Vehtari A (2019). *R-squared for Bayesian Regression Models*. *The American Statistician* **73**, 307–309.
- 665
- 666 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.
- 667
- 668 Gilmour AR, Thompson R, Cullis BR (1995). *Average Information REML: An Efficient Algorithm for Variance Parameter Estimation in Linear Mixed Models*. *Biometrics*, 1440–1450.
- 669
- 670 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, de 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>.
- 671
- 672
- 673
- 674
- 675 Hampel FR, Ronchetti EM, Rousseeuw PJ, Stahel WA (2011). *Robust Statistics: The Approach Based on Influence Functions*. John Wiley & Sons.
- 676
- 677 Helland IS (1987). *On the interpretation and use of R^2 in regression analysis*. *Biometrics*, 61–69.
- 678
- 679 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.
- 680
- 681 Huber PJ, Ronchetti EM (1981). *Robust Statistics*, Wiley: New York. USA.
- 682
- 683 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>.
- 684
- 685 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>.
- 686
- 687
- 688 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.
- 689
- 690
- 691 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.
- 692
- 693
- 694 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.
- 695
- 696
- 697 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.
- 698
- 699
- 700 Lange KL, Little RJ, Taylor JM (1989). *Robust Statistical Modeling Using the t Distribution*. *Journal of the American Statistical Association* **84**, 881–896.
- 701
- 702 Lartillot N (2023). *Identifying the best approximating model in Bayesian phylogenetics: Bayes factors, cross-validation or wAIC?* *Systematic Biology* **72**, 616–638.
- 703
- 704 Lian L, de los 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>.
- 705
- 706
- 707 Lin CS, Binns MR, Lefkovitch LP (1986). *Stability Analysis: Where Do We Stand?* *Crop science* **26**, 894–900.
- 708
- 709 Martin A, Sherington J (1997). *Participatory Research Methods—Implementation, Effectiveness and Institutional Context*. *Agricultural systems* **55**, 195–216.
- 710

- 711 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.
- 712
- 713
- 714 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>.
- 715
- 716
- 717 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.
- 718
- 719 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>. JSTOR: 1400388.
- 720
- 721
- 722 Nascimento M, Nascimento ACC, e Silva FF, Teodoro PE, Azevedo CF, de Oliveira TRA, do Amaral Junior AT, Cruz CD, Farias FJC, de Carvalho LP (2020). *Bayesian Segmented Regression Model for Adaptability and Stability Evaluation of Cotton Genotypes*. *Euphytica* **216**, 30. <https://doi.org/10.1007/s10681-020-2564-5>.
- 723
- 724
- 725
- 726 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.
- 727
- 728
- 729 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>.
- 730
- 731 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.
- 732
- 733 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.
- 734
- 735 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.
- 736
- 737
- 738 Pereira DG, Mexia JT, Rodrigues PC (2007). *Robustness of Joint Regression Analysis*. *Biometrical Letters* **44**, 105–128.
- 739
- 740 Perkins JM, Jinks JL (1968). *Environmental and Genotype-Environmental Components of Variability*. *Heredity* **23**, 339–356.
- 741
- 742 Piepho HP (1999). *Stability analysis using the SAS system*. *Agronomy Journal* **91**, 154–160.
- 743
- 744 Piepho HP (2004). *An algorithm for a letter-based representation of all-pairwise comparisons*. *Journal of Computational and Graphical Statistics* **13**, 456–466.
- 745
- 746 Piepho HP, Blancon J (2023). *Extending Finlay–Wilkinson regression with environmental covariates*. *Plant Breeding* **142**, 621–631.
- 747
- 748 R Core Team (2014). *R: A Language and Environment for Statistical Computing*.
- 749
- 750 Reckling M, Ahrends H, Chen TW, Eugster W, Hadasch S, Knapp S, Laidig F, Linstädter A, Marcholdt J, Piepho HP (2021). *Methods of Yield Stability Analysis in Long-Term Field Experiments. A Review*. *Agronomy for Sustainable Development* **41**, 1–28.
- 751
- 752 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>.
- 753
- 754 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>.
- 755
- 756
- 757 Robert C (2007). *The Bayesian Choice: From Decision-Theoretic Foundations to Computational Implementation*. Springer Science & Business Media.
- 758
- 759 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.
- 760
- 761
- 762 Rosa Padovani CR, Gianola D (2003). *Robust Linear Mixed Models with Normal/Independent Distributions and Bayesian MCMC Implementation*. *Biometrical Journal* **45**, 573–590.
- 763
- 764 Silvey SD (1975). *Statistical Inference*. Chapman and Hall.

- 765 Simar L (2002). *Les modèles de base de l'analyse bayésienne*. In: *Méthodes Bayésiennes en statistique*.
766 Ed. by Jean-Jacques Dreesbeke, Jeanne Fine, and Gilbert Saporta. Editions Technip, pp. 103–
767 132.
- 768 Smith Cullis BR, Thompson R (2005). *The Analysis of Crop Cultivar Breeding and Evaluation Trials:*
769 *An Overview of Current Mixed Model Approaches*. *The Journal of Agricultural Science* **143**, 449–
770 462.
- 771 Smith ME, Castillo FG, Gómez F (2001). *Participatory Plant Breeding with Maize in Mexico and*
772 *Honduras*. *Euphytica* **122**, 551–563. <https://doi.org/10.1023/A:1017510529440>.
- 773 Snapp SS, Silim SN (2002). *Farmer Preferences and Legume Intensification for Low Nutrient Environ-*
774 *ments*. *Plant and soil* **245**, 181–192.
- 775 Stan Development Team (2016). *RStan: The R Interface to Stan*.
- 776 Talbot M (1984). *Yield variability of crop varieties in the UK*. *The Journal of Agricultural Science* **102**,
777 315–321.
- 778 Thompson R, Cullis B, Smith A, Gilmour A (2003). *A Sparse Implementation of the Average Informa-*
779 *tion Algorithm for Factor Analytic and Reduced Rank Variance Models*. *Australian & New Zealand*
780 *Journal of Statistics* **45**, 445–459.
- 781 Turbet Delof M (2024). *Impacts de l'environnement sur les pratiques de sélection paysanne et le*
782 *comportement des variétés qui en résultent*. PhD thesis. Université Paris-Saclay.
- 783 van Etten J, Sousa K, Aguilar A, Barrios M, Coto A, Dell'Acqua M, Fadda C, Gebrehawaryat Y,
784 Gevel J, Gupta A, et al. (2019). *Crop variety management for climate adaptation supported by*
785 *citizen science*. *Proceedings of the National Academy of Sciences* **116**, 4194–4199.
- 786 van Frank G (2018). *Gestion participative de la diversité cultivée et création de mélanges diversifiés*
787 *de blé tendre à la ferme*. PhD thesis. Université Paris Saclay (COMUE).
- 788 van Frank G, Goldringer I, Rivière P, David O (2019). *Influence of Experimental Design on Decen-*
789 *tralized, on-Farm Evaluation of Populations: A Simulation Study*. *Euphytica* **215**, 126. <https://doi.org/10.1007/s10681-019-2447-9>.
- 790
- 791 van Frank G, Rivière P, Pin S, Baltassat R, Berthelot JF, Caizergues F, Dalmasso C, Gascuel JS,
792 Hyacinthe A, Mercier F, Montaz H, Ronot B, Goldringer I (2020). *Genetic Diversity and Stability*
793 *of Performance of Wheat Population Varieties Developed by Participatory Breeding*. *Sustainability*
794 **12**, 384. <https://doi.org/10.3390/su12010384>.
- 795 Vehtari A, Gelman A, Gabry J (2017). *Practical Bayesian Model Evaluation Using Leave-One-out*
796 *Cross-Validation and WAIC*. *Statistics and Computing* **27**, 1413–1432. [https://doi.org/10.](https://doi.org/10.1007/s11222-016-9696-4)
797 [1007/s11222-016-9696-4](https://doi.org/10.1007/s11222-016-9696-4).
- 798 Virk DS, Chakraborty M, Ghosh J, Prasad SC, Witcombe JR (2005). *Increasing the Client Orienta-*
799 *tion of Maize Breeding Using Farmer Participation in Eastern India*. *Experimental Agriculture* **41**,
800 413–426. <https://doi.org/10.1017/S001447970500270X>.
- 801 Walter E, Pronzato L (1997). *Identification of parametric models: from experimental data*. Springer.
- 802 Welham SJ, Gogel BJ, Smith AB, Thompson R, Cullis BR (2010). *A comparison of analysis methods*
803 *for late-stage variety evaluation trials*. *Australian & New Zealand Journal of Statistics* **52**, 125–
804 149.
- 805 Wolfe MS, Baresel JP, Desclaux D, Goldringer I, Hoad S, Kovacs G, Löschenberger F, Miedaner T,
806 Østergård H, Lammerts van Bueren ET (2008). *Developments in Breeding Cereals for Organic*
807 *Agriculture*. *Euphytica* **163**, 323. <https://doi.org/10.1007/s10681-008-9690-9>.
- 808 Woyann LG, Benin G, Storck L, Trevizan DM, Meneguzzi C, Marchioro VS, Tonatto M, Madureira
809 A (2017). *Estimation of Missing Values Affects Important Aspects of GGE Biplot Analysis*. *Crop*
810 *Science* **57**, 40–52.
- 811 Wricke G (1962). *Über Eine Methode Zur Erfassung Der Okologischen Streubreite in Feldversuchen*.
812 *Z. Pflanzenzuchtg* **47**, 92–96.
- 813 Yan W (2013). *Biplot Analysis of Incomplete Two-way Data*. *Crop Science* **53**, 48–57.
- 814 Yates F, Cochran WG (1938). *The Analysis of Groups of Experiments*. *The Journal of Agricultural*
815 *Science* **28**, 556–580. <https://doi.org/10.1017/S0021859600050978>.