

1 FW: An R package for Finlay-Wilkinson Regression that Incorporates
2 Genomic/Pedigree Information and Covariance Structures Between Environments

3 Lian Lian^{1*} and Gustavo de los Campos^{1,2}

4 1. Department of Epidemiology & Biostatistics, Michigan State University, 909 Fee
5 Road, Room B601, East Lansing, MI, 48824.
6 2. Department of Probability and Statistics, Michigan State University.
7 *Corresponding author (lianl0501@gmail.com)

8 **Abstract**

9 The Finlay-Wilkinson Regression is a popular method among plant breeders to
10 describe genotype by environment interaction. The standard implementation is a two-step
11 procedure that uses environment (sample) means as covariates in a within-line ordinary
12 least squares (OLS) regression. This procedure can be suboptimal for at least four
13 reasons: (i) in the first step environmental means are typically **estimated** without
14 **considering genetic-by-environment interactions**, (ii) in the second step uncertainty about
15 the environmental means is ignored, (iii) estimation is performed regarding lines and
16 environment as fixed effects and (iv) the procedure does not incorporate genetic (either
17 pedigree-derived or marker derived) relationships. Su et al. proposed to address these
18 problems using a Bayesian method that allows simultaneous estimation of environmental
19 and genotype parameters, and allows incorporation of pedigree information. In this article
20 we: (i) extend the model presented by Su et al. to allow integration of **genomic (e.g.,**
21 **SNP)** and **covariance between environments**, (ii) present an R package (FW) that
22 implements these methods, and (iii) illustrate the use of the package using examples
23 based on real data. The FW R-package implements both the two-step OLS method and a
24 full Bayesian approach for Finlay-Wilkinson regression with a very simple interface.

1 Using a real wheat data set we demonstrate that the prediction accuracy of the Bayesian
2 approach is consistently higher than the one achieved by the two-steps OLS method.

3 **Introduction**

4 Plant breeders use the Finlay-Wilkinson Regression (FW, Finlay and Wilkinson,
5 1963) to assess stability of varieties across different environments. The FW aims at
6 assessing how the expected performance of a genotype varies as a function of the
7 environmental effects. Usually this is achieved by regressing the performance of each
8 genotype on the environmental means. Compared with a completely un-structured
9 genotype by environment interaction ($G \times E$) model that fits every level of genotype and
10 environment combination, the Finlay-Wilkinson regression is parsimonious and can
11 reveal a trend of variety performance across environments. Breeders can use this model
12 to select for plants either based on stability or on responsiveness to environment potential
13 (Walsh and Lynch, 2014).

14 The standard implementation of Finlay-Wilkinson regression is a two-step
15 procedure whereas in the first step environmental sample means are computed and in the
16 second step intercepts and slopes of each line are estimated by regressing, within line, the
17 performance of each line on the estimated environmental means. This procedure has at
18 least four potential limitations: (i) in the first step environmental means are typically
19 estimated without considering $G \times E$, (ii) in the second step, uncertainty about the
20 environmental means is ignored (iii) the environmental means and the variety intercepts
21 and slopes are regarded as fixed effects (this can lead to large sampling variance of
22 estimates), and (iv) the procedure does not offer a clear way of incorporating pedigree or

1 molecular marker information when estimating the intercepts and slopes of the lines.
2 These drawbacks can induce biases (especially in incomplete designs where a few lines
3 are evaluated in each environments) and lead to large sampling variance of estimates.

4 Su et al. (2006) proposed a Bayesian method that addresses the limitations of the
5 standard two-step procedure. The methodology described by Su et al.: (1) uses a Gibbs
6 sampler that allows estimating environmental and genotype parameters jointly, (2) fully
7 accounts for confounding and uncertainty about environmental means, (3) treats
8 environmental means and the intercepts and slopes of the lines as random—this treatment
9 usually perform better than ordinary least squares in terms of mean-squared error and of
10 prediction accuracy, especially when the number of parameters to be estimated is large
11 relative to sample size (Copas, 1983; Frank and Friedman, 1993), and (iv) allows
12 incorporating pedigree information into the model. Using simulations, Su et al. (2006)
13 reported better statistical performance of the Bayesian method for estimating model
14 parameters. In this article we extend the model proposed by Su et al. (2006) in ways that
15 allow incorporating genomic (e.g., SNP) information and covariance between the
16 environment effects.

17 To the best of our knowledge the methodology described by Su et al. for animal
18 breeding applications has not been considered in plant breeding, and there is no publicly
19 available user-friendly software for implementing a Bayesian FW regression. Therefore,
20 in this article we introduce an R-package (R Development Core Team, 2011) that
21 implements the Finlay Wilkinson regression. The FW package implements both the two-
22 steps ordinary least squares (OLS) procedure and Bayesian single step procedure that
23 allows incorporating covariance structure for varieties (e.g., a pedigree or marker-derived

1 kinship matrix) and environments. We describe the methods implemented in the package
2 and show with examples how this package can be used to perform the Finlay-Wilkinson
3 regression with both methods. Finally, we present an evaluation of prediction accuracy
4 for the Bayesian and two-step OLS methods with a wheat data set.

5 **Model Specification and Algorithm**

6 In a reaction norm model (Gregorius and Namkoong, 1986; Perkins and Jinks,
7 1968) the phenotypic record of the k th replicate of the i th variety observed in the j th
8 environment is modeled as follows

$$9 \quad y_{ijk} = \mu + g_i + h_j + b_i h_j + \varepsilon_{ijk} \quad [\text{Eq. 1}]$$

10 where g_i is the main effect of i^{th} variety and h_j is the main effect of the j^{th} environment,
11 and ε_{ijk} is an error term, usually assumed to be IID normal with mean zero and variance
12 σ_ε^2 . When we reorganize Eq. 1 into the form: $y_{ijk} = \mu + g_i + (b_i + 1)h_j + \varepsilon_{ijk}$, we can
13 recognize that $b_i + 1$ is the change of expected variety performance per unit change of
14 the environment effect (h_j). If there are no replicates the index k can be removed. With
15 this, the equation reduces to $y_{ij} = \mu + g_i + h_j + b_i h_j + \varepsilon_{ij}$. The collection of
16 parameters to be estimated from the model of Eq. 1 include the intercept and the vectors
17 of effects: $\mathbf{g} = \{g_i\}$, $\mathbf{b} = \{b_i\}$ and $\mathbf{h} = \{h_j\}$.

18 **Estimation using two steps methods**

19 The Finlay-Wilkinson regression requires regressing the observed phenotypes of
20 the line on environment effects. In the standard Finlay-Wilkinson regression (Finlay and
21 Wilkinson, 1963) the environmental effects are computed from the sample environmental
22 means. However, in incomplete designs the sample mean of an environment may not be

an un-biased estimate of the true environment mean. Therefore, a better estimate of environment effects comes from a regression that accounts for both environment effects and genotype effects, that is:

Step 1, estimate the environmental effect using a main effects model

$$y_{ijk} = \mu + g_i + h_j + \varepsilon_{ijk} \quad [\text{Eq. 2}]$$

The above regression yields estimates of environment effects (\hat{h}_j), these can be used in the second step to estimate the intercepts and slopes of each line.

Step 2, replace h_j with \hat{h}_j in equation [1] yielding

$$y_{ijk} = \mu + g_i + \hat{h}_j + b_i \hat{h}_j + \varepsilon_{ijk} \quad [\text{Eq. 3}]$$

Both Eq.2 and Eq.3 can be implemented with either ordinary least squares (OLS) or mixed models. The current FW package implemented both Step 1 and Step 2 with OLS. In Step 1, Eq.2 is fitted with the constraint that $\sum_j \hat{h}_j = 0$ and $\sum_i \hat{g}_i = 0$. Step 2 is implemented by fitting Eq. 3 separately within each line with the constraint $\hat{\mu} = 0$.

Bayesian approach

Bayesian inferences are based on the posterior distribution of unknown parameters given the data: $p(\boldsymbol{\theta}|\mathbf{y}) \propto p(\mathbf{y}|\boldsymbol{\theta})p(\boldsymbol{\theta})$, where $\boldsymbol{\theta}$ represents the collection of the unknowns: $\boldsymbol{\theta} = \{\mu, \mathbf{g}, \mathbf{b}, \mathbf{h}, \sigma_g^2, \sigma_b^2, \sigma_h^2, \sigma_\varepsilon^2\}$, $p(\mathbf{y}|\boldsymbol{\theta})$ is the conditional distribution of the data given the parameters and $p(\boldsymbol{\theta})$ is the joint prior distribution assigned to the model unknowns. According to Eq. 1 and assuming IID normal residuals, we have

$$p(\mathbf{y}|\boldsymbol{\theta}) = \prod_{ijk} N(\mu + g_i + h_j + b_i h_j, \sigma_\varepsilon^2).$$

In the FW package, the prior density is assumed to have the following form:

$$p(\boldsymbol{\theta}) = p(\sigma_\varepsilon^2)p(\mathbf{g}|\sigma_g^2)p(\mathbf{b}|\sigma_b^2)p(\mathbf{h}|\sigma_h^2)p(\sigma_g^2)p(\sigma_b^2)p(\sigma_h^2). \text{ The residual variance } \sigma_\varepsilon^2 \text{ is}$$

1 assigned a scaled-inverse χ^2 distribution: $\sigma_\varepsilon^2 \sim \chi^{-2}(\nu_\varepsilon, S_\varepsilon^2)$, with degrees of freedom ν_ε
 2 (>0) and scale parameter S_ε^2 (>0), in the parameterization used $E[\sigma_\varepsilon^2] = \frac{\nu_\varepsilon S_\varepsilon^2}{\nu_\varepsilon - 2}$. The
 3 overall mean μ is assigned a flat prior. The prior distributions for \mathbf{g} , \mathbf{b} , \mathbf{h} , are all
 4 multivariate Normal: $\mathbf{h} \sim N(\mathbf{0}, \mathbf{H}\sigma_h^2)$, $\mathbf{g} \sim N(\mathbf{0}, \mathbf{A}\sigma_g^2)$, $\mathbf{b} \sim N(\mathbf{0}, \mathbf{A}\sigma_b^2)$, where \mathbf{H} is a
 5 covariance structure describing co-variances between the environmental means (this can
 6 be a covariance structure based on spatial information) and \mathbf{A} is a covariance structure
 7 describing co-variances between levels of the random effects \mathbf{g} and \mathbf{b} (\mathbf{A} could be either a
 8 pedigree or marker-derived relationship matrix). Independence between the effects of the
 9 levels of any of the random effects can be obtained by setting either \mathbf{A} or \mathbf{H} to be an
 10 identity matrix. The variance components σ_h^2 , σ_g^2 and σ_b^2 are assigned scaled-inverse- χ^2
 11 distributions whose shape are controlled by variance-specific degree of freedom and scale
 12 hyper parameters. The FW package offers users the possibility of specifying hyper
 13 parameters (degree of freedom and scale parameters); however, if these are not specified,
 14 specific sets of rules similar to those described in (Pérez et al., 2010) are used to
 15 determine those parameters. Further details about this are given in the supplemental files.

16 In the model described above the posterior density does not have a closed form;
 17 however, estimates of features of the posterior distribution (e.g., posterior means,
 18 posterior standard deviations, or credibility regions) can be derived using Monte Carlo
 19 Methods. The FW package draws samples from the posterior distribution of the model
 20 using a Gibbs sampler (Casella and George, 1992; Geman and Geman, 1984) similar to
 21 the one described in (Su et al., 2006), details of the implementation of Gibbs sampler are
 22 provided in the supplemental files.

1 **Software**

2 The FW package implements both a two-steps OLS method and the Bayesian
3 method described in the previous section. Type the following command in R will install
4 the package:

```
5        library(devtools)  
6        install_github("lian0090/FW")
```

8 **Wheat data set**

9 The package includes a data set that can be used to run examples. The data set
10 (originally made publicly available by Crossa et al., 2010) contains data for 599 wheat
11 lines from CIMMYT's Global Wheat Program and evaluated for grain yield in four
12 environments. The dataset becomes available in the R environment by running the
13 following R-code:

```
14        library(FW)  
15        data(wheat)
```

17 Function `library()` loads the package, and `data()` loads datasets included
18 in the package into the environment. The above code loads the following objects into the
19 environment: (i) `wheat.Y`, a `data.frame` (2396×3) containing the grain yield (average
20 of two plot records, $\$Y$) of 599 wheat lines ($\VAR) in four environments ($\$ENV$), (ii)
21 `wheat.G` (599×599) is a genomic relationship matrix computed from DArT markers.
22 Further details about this data set can be found in Crossa et al. (2010).

1 User Interface

2 All the arguments of the FW function have default values, except the response
3 variable and the corresponding identifiers for varieties and environments. A basic call to
4 the FW program is as follows.

Box 1: Basic Call of the FW Function

```
1 library(FW)
2 data(wheat)
3 attach(wheat.Y)
4 lm1=FW(y=y, VAR=VAR, ENV=ENV, method="OLS")
5 lm2=FW(y=y, VAR=VAR, ENV=ENV)
```

5

6 When the call of the FW function is done using the code in line 4 of Box 1, FW
7 fits a Finlay-Wilkinson regression with the two-steps OLS method: *y* (numeric, *n*, NAs
8 are allowed) is the response variable, *VAR* (character, *n*, NAs are not allowed) are the
9 identifiers for the varieties which are treated as labels; *ENV* (character, *n*, NAs are not
10 allowed) are the identifiers for the environments; *method* is used to describe what
11 method to use: "OLS" for ordinary least squares. The default method ("Gibbs") is the
12 Bayesian regression; this can be invoked using the code in line 5 of Box 1. By default, a
13 single chain of Gibbs sampler is run with a total of 5,000 cycles and the samples from the
14 first 3,000 cycles are used for Burn-in, samples from the remaining 2,000 cycles for
15 inference (the user is advised to run longer chains and to check convergence as well as
16 the size of Monte Carlo errors). The FW function provides many additional arguments
17 that can be used to specify the model (e.g., providing co-variance matrices for varieties
18 and environments, user-defined values for hyper-parameters) and the algorithm (number

1 of chains, numbers of iterations, etc.); details can be found in the user manual and in the
2 examples presented below.

3 After fitting either OLS or Gibbs method, FW function returns a list with
4 estimates and arguments used in the call, a brief description of the outputs follows.

5 **Return**

6 Box 2 shows the structure of the object returned after calling the FW function (see
7 line 1 of Box 2). The first element `$y` of the list is the response vector used in the call to
8 FW, `$whichNa` gives the position of the entries in `y` that were missing, `$mu` (vector),
9 `$g` (matrix), `$b` (matrix), `$h` (matrix) are the estimated posterior means of μ , \mathbf{g} , \mathbf{b} , \mathbf{h} ;
10 `$yhat` (matrix) is the estimated posterior means of the predictor $\hat{\mathbf{y}}$: $\hat{y}_{ijk} = \hat{\mu} + \hat{g}_i + \hat{h}_j +$
11 $\hat{b}_i \hat{h}_j$; `$SD.mu` (vector), `$SD.g` (matrix), `$SD.b` (matrix), `$SD.h` (matrix) and
12 `$SD.yhat` are the estimated posterior standard deviation for μ , \mathbf{g} , \mathbf{b} , \mathbf{h} and $\mu + g_i +$
13 $h_j + b_i h_j$ respectively.

14 With OLS method, `$g`, `$b`, `$h` and `$yhat` all have only one column; with Gibbs
15 method each column provides estimates derived from one MCMC chain. Since the
16 default behavior is to run only one chain the outputs in Box 2 contain only one column;
17 however, if multiple chains are run, estimates from different chains are provided in
18 different columns.

19 The output `$var_e`, `$var_g`, `$var_b`, `$var_h` are the estimated posterior
20 means for σ_ε^2 , σ_g^2 , σ_b^2 and σ_h^2 (only available for Gibbs method). Each element of
21 `$var_e`, `$var_g`, `$var_b` and `$var_h` correspond to estimates derived from

1 different chains; `$SD.var_e`, `$SD.var_g`, `$SD.var_b` and `$SD.var_h` are the
2 estimated posterior standard deviation for σ_ε^2 , σ_g^2 , σ_b^2 and σ_h^2 , respectively.

3

Box 2: Structure of the object returned by FW

```

1 str(lm2)
2 List of 24
3  $ y      : num [1:2396] 6.17 3.14 2.74 3.26 4.99 ...
4  $ whichNa : int(0)
5  $ VAR     : chr [1:2396] "775" "775" "775" "775" ...
6  $ ENV     : chr [1:2396] "1" "2" "4" "5" ...
7  $ mu      : Named num 4.64
8  $ SD.mu   : Named num 0.0979
9  $ g       : num [1:599, 1] -0.476 0.16 -0.611 ...
10 $ SD.g     : num [1:599, 1] 0.224 0.219 0.224 0.208 ...
11 $ b       : num [1:599, 1] 0.1604 -0.1255 0.251 ...
12 $ SD.b     : num [1:599, 1] 0.237 0.236 0.235 0.24 ...
13 $ h       : num [1:4, 1] 0.519 -0.186 -0.776 -1.383 ...
14 $ SD.h     : num [1:4, 1] 0.096 0.0999 0.0999 0.103 ...
15 $ yhat     : num [1:2396, 1] 5.17 4.3 3.56 2.81 5.21 ...
16 $ SD.yhat  : num [1:2396, 1] 0.283 0.217 0.25 0.343 ...
17 $ var_e    : Named num 0.3
18 $ SD.var_e : Named num 0.0111
19 $ var_g    : Named num 0.0885
20 $ SD.var_g : Named num 0.0116
21 $ var_b    : Named num 0.0973
22 $ SD.var_b : Named num 0.0132
23 $ var_h    : Named num 0.926
24 $ SD.var_h : Named num 0.595

```

4 Output files

5 No output files are generated for OLS method. For Gibbs method, samples for σ_ε^2 ,
6 σ_g^2 , σ_b^2 , σ_h^2 , and (by default) the first two elements of **g**, **b**, **h** will be saved; as the Gibbs
7 sampler collects samples, these samples are saved to the hard drive (only the most recent
8 samples are retained in memory); by default, a thinning of 5 is used. Once the iteration
9 process finishes, FW will read all the saved samples into a `mcmc` object, save the `mcmc`
10 object into a file `samps.rda`, and remove the raw sample files. To prevent overloading
11 the RAM with samples by default FW only save samples of the first two entries of the
12 vectors of random effects; however the user can change this behavior by specifying

1 which entries of the vectors are desired using the `saveVAR` (for **g** and **b**) and `saveENV`
2 (for **h**) argument. These samples produced by FW can be used to assess convergence and
3 to estimate Monte Carlo Standard Errors. The file `samps.rda` can be directly loaded
4 into R using `load('samps.rda')`. Once the object containing the samples is loaded
5 in the R environment, the package `coda` (Plummer et al., 2006) can be used to obtain
6 plots of the chains and compute convergence diagnostics.

7 **Application examples**

8 In this section we illustrate via examples some of the features of the FW package;
9 Example 1 illustrates how the package can be used to fit Finlay-Wilkinson regression by
10 OLS method and Gibbs method with and without covariance structure and Example 2
11 describes how the package can be used for cross-validation analyses. Additional
12 examples involving fine-tuning the Gibbs method (e.g., hyper-parameter setup, fitting
13 more than two chains, specify saved samples) are provided as Supplementary data.

14 **Example 1: Fitting models with default setup for 599 wheat lines**

15 Box 3 shows the code used to fit a FW regression using three different
16 approaches: (i) a two-steps OLS method (code in line 3), (ii) a Bayesian FW regression
17 assuming independence of lines and of environments (code in lines 5-7) and (iii) a
18 Bayesian FW regression that incorporates genomic information (lines 9-11). In the
19 Bayesian models, the seed for the random number generator can be specified using the
20 argument `seed` (see lines 5-11) and the argument `saveAt` can be used to add a path and
21 a pre-fix to be appended to `'samps.rda'` file.

22

1

Box 3: Fit models by default parameters

```

1 library(FW); data(wheat); attach(wheat.Y)
2
3 OLS=FW(y=y, VAR=VAR, ENV=ENV, method="OLS")
4
5 GibbsI=FW( y=y, VAR=VAR, ENV=ENV,
6           method="Gibbs", seed=12345, saveAt="GibbsI", nIter=50000
7           , burnIn=5000)
8
9 GibbsA=FW(y=y, VAR=VAR, ENV=ENV,
10          method="Gibbs", A=wheat.G, seed=12345,
11          saveAt="GibbsA", nIter=50000, burnIn=5000)
12
13 load("GibbsIsamps.rda")
14 HPDinterval(samps[,c("var_e", "var_g", "var_b", "var_h")])
15
16 load("GibbsAsamps.rda")
17 HPDinterval(samps[,c("var_e", "var_g", "var_b", "var_h")])

```

2

3 Parameter estimates (estimated posterior means) can be directly extracted from
4 the FW object as illustrated in Box 2. Other features of the posterior distribution (e.g., 95%
5 credibility intervals for the parameters) can be obtained by post-hoc analyses of the
6 samples included in the `rda` file generated by the program (see, line 13-17 of Box 3). In
7 Table 1, we listed the estimates of variance components from the three models. For OLS
8 method, only the residual variance σ_ϵ^2 (weighted mean of residual variance for each
9 within line regression by its residual degree of freedom) is estimated. The estimated error
10 variances are very similar across the three models. Also from Table 1, we can see that the
11 estimated variance of the main effects of the environments is large relative to both the
12 error variance and the phenotypic variance.

13

14 **Table 1.** Estimated variance components (posterior 95% credibility intervals in
15 parenthesis) from different models

Parameters	FW output	OLS	GibbsI (A=I)	GibbsA (A=G)
σ_{ε}^2	\$var_e (Gibbs) \$var_e_weighted(OLS)	0.32	0.30 (0.28, 0.32)	0.30 (0.28, 0.32)
σ_g^2	\$var_g	NA	0.09 (0.07, 0.11)	0.11 (0.08, 0.14)
σ_b^2	\$var_b	NA	0.10 (0.07, 0.12)	0.13 (0.10, 0.17)
σ_h^2	\$var_h	NA	0.90 (0.24, 1.90)	0.88 (0.24, 1.88)

The fitness of the models can be examined by the correlations between the observed values y and the fitted values \hat{y} (line 1 of Box 4). The OLS model fitted the data better than both GibbsI and GibbsA: the correlation was 0.91 for the OLS method, 0.88 those for GibbsI and 0.86 for GibbsA.

Box 4: Correlation between y and \hat{y} , and correlations for \hat{b} among different models.

```

1 cor(y, OLS$yhat); cor(y, GibbsI$yhat); cor(y, GibbsA$yhat);
2 cor(OLS$b, GibbsI$b); cor(OLS$b, GibbsA$b); cor(GibbsI$b, GibbsA$b)

```

In Table 2, we listed the correlations among parameter estimates from different models (code for \hat{b} was provided in line 2 of Box 4), and noticed that correlations among parameters estimates from different models are high; this is expected considering that the data comes from a full factorial design where all lines are evaluated in all environments.

Table 2. Pearson's product-moment correlation between parameter estimates derived by each of the three methods implemented in Box 3.

OLS-GibbsI	OLS-GibbsA	GibbsI-GibbsA
------------	------------	---------------

\hat{h}	1.00	1.00	1.00
\hat{b}	0.94	0.81	0.83
\hat{g}	0.98	0.79	0.81
\hat{y}	0.96	0.94	0.97

1

2 The pattern of variety performance in different environments can be visualized by
3 plotting the observed and fitted values against the estimated environment effects. Figure
4 1 was generated by the calling of `plot` function in line 2-3 of Box 5. Each line in this
5 plot corresponds to a genotype. The comparison of the results from the OLS and GibbsA
6 reveals interesting patterns: the OLS method predicts a much stronger extent of
7 variability in intercepts and slopes (this is likely due to over-fitting, see Example 2
8 below) than the Bayesian method. The Bayesian method yields ‘smoother’ predictions;
9 this is a direct consequence of the shrinkage-towards-the-mean induced in the Bayesian
10 method by treating effects as random and the use of correlations between genotypes (e.g.,
11 genomic relationships).

12 The function `plotVAR` also allows users to display the curves for a few
13 genotypes (see code in lines 6-11 of Box5). Using this feature we display in Figure 2 the
14 estimated regressions for five varieties. The slope in the plot corresponds to $1 + b_i$ and
15 the dashed gray line corresponds to a slope equals to 1 ($b_i = 0$), recall that $1 + b_i$
16 represents the expected change in performance of the i^{th} variety per unit change in the
17 environment effect. We observe from Figure 2 that line ID=1081265 performs well in all
18 environments and line ID=13302 is better adapted to good environments.

19

Box 5: Plot fitted models

```

1 par(mfrow=c(1,2))
2 plot(OLS,main="OLS", cex=0.2,lwd=0.2)
3 plot(GibbsA,main="GibbsA", cex=0.2,lwd=0.2) #cex controls point
4 size, lwd controls the line width
5
6 plot(OLS, plotVAR=c("1081265","1101307",
7 "1295736", "13302" , "1343502"), main="OLS")
8
9 plot(GibbsA, plotVAR=c("1081265","1101307",
10 "1295736", "13302" , "1343502"),
11 main="GibbsA")

```

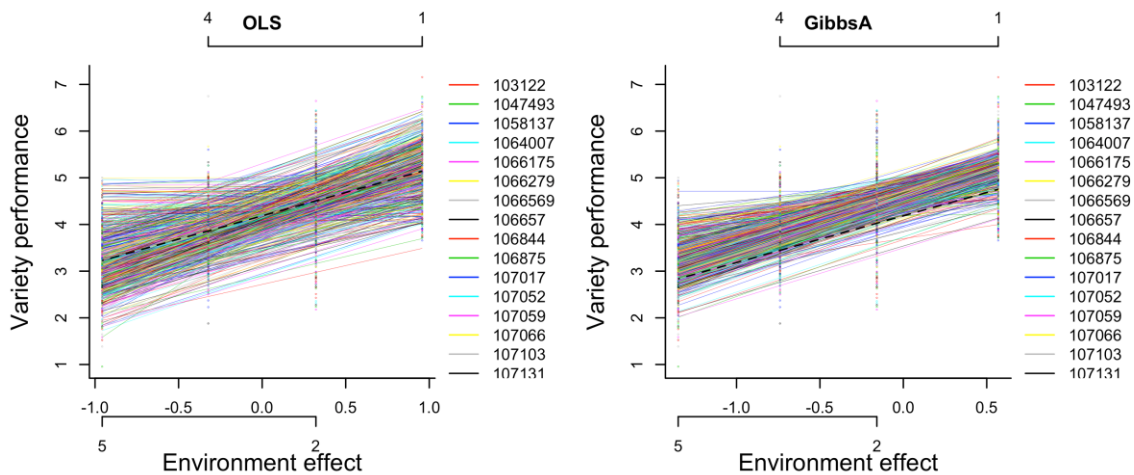


Figure 1. Plot of variety performance versus estimated environment values. Each line represents a different variety. Lines are fitted values and points are the cell means of genotype and environment combination. The horizontal axis in Figure 1 displays the estimated environmental effects. The labels of these environments are also displayed, these labels can be removed by setting `ENVlabel=F`.

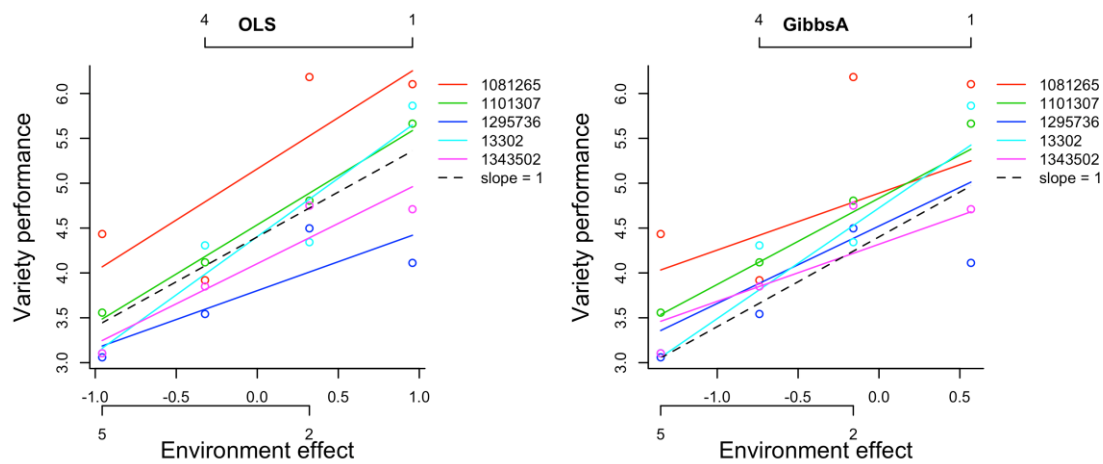


Figure 2. Plot of the performance of five varieties on estimated environment values.

Each color represents a different variety. Lines are fitted values and circles are the cell means of genotype by environment combination. The horizontal axis in Figure 2 displays the estimated environmental effects. The labels of these environments are also displayed, these labels can be removed by setting `ENVlabel=F`.

Fitting models with covariance between the environment effects

Covariance structures can be used to induce borrowing of information between levels of a random effect. For instance, pedigree-based or genomic-derived relationships can be used to induce borrowing of information between genotypes. Similarly a covariance structure between the environment effects could be used to induce borrowing of information between environments. Such covariance structures can be derived from previous knowledge about the correlation of the average performance of genotypes in pairs of environments or by using environmental covariates, as demonstrated in (Jarquín

et al., 2014). The FW package allows incorporating covariance between the environment effects; an example of how this can be done is given in Box 6. In this example we compare three analyses. The first model (`GibbsI`) assumes that the environment effects are independent; this model was fitted previously using the code in Box 3. Subsequently, we modified this model by incorporating a covariance structure that assumes a covariance of 0.9 between environments 1 and 2, and null covariance among the other pairs of environments. This model was fitted using the entire data set (`GibbsH`) and after setting to NA all the records from the 2nd environment (`GibbsH_NA`).

Table 3 displays the estimated environment effects derived from each of the analyses. The estimated environment effects derived from `GibbsI` and `GibbsH` were almost identical. This happens because in these two examples the data available for each environment dominates over the prior (which in case of `GibbsH` assumes that the effects of environments 1 and 2 are highly correlated). However, when we set to NA all the entries of environment 2 (`GibbsH_NA`), the estimated effects for environments 1 and 2 are very close. This was entirely driven by the covariance structure H. An intermediate situation can emerge where one environment has records for a few genotypes. In such cases, non-diagonal covariance structures (H) may be used to borrow information between environments.

Finally, the example provided by `GibbsH_NA` also illustrates how H allows to make predictions about environments without records; if such environments are correlated with other environments for which we have data, in principle we can infer the effects for those environments. This of course won't be possible if H is diagonal.

1 Table 3 Estimated environment effects from GibbsI and GibbsH

ENV	GibbsI	GibbsH	GibbsH2
1	0.52	0.51	0.78
2	-0.18	-0.19	0.74
4	-0.78	-0.78	-0.52
5	-1.38	-1.39	-1.11

2

Box 6: Including covariance matrix (H) for Environments in FW

```

1 H=diag(1, 4)
2 H[1, 2]=H[2, 1]=0.9
3 colnames(H)=rownames(H)=c(1, 2, 4, 5)
4
5 GibbsH=FW(y=y, VAR=VAR, ENV=ENV,
6           method="Gibbs", H=H, seed=12345, nIter=50000, burnIn=5000)
7
8 yNA=y
9 yNA[which(ENV==2)]=NA
10
11 GibbsH_NA=FW(y=yNA, VAR=VAR, ENV=ENV,
12              method="Gibbs", H=H, seed=12345, nIter=50000, burnIn=5000)
13
14 round(cbind(GibbsI$h, GibbsH$h, GibbsH_NA$h), 2)

```

3 Assessment of convergence for Bayesian FW regressions

4 The convergence of Gibbs sampler can be examined by plotting the samples
5 collected by FW. The code in Box 7 illustrates how to produce trace plots: lines 1-2 load
6 and plot the samples from GibbsI and lines 3-4 do the same for GibbsA. Mixing was
7 reasonably good (samples traverse through the sample space in relatively few steps, and
8 can be verified by low average autocorrelation between samples: for example the average
9 autocorrelation was 0.05 for var_e at lag 5, see line 5 of Box 7) in both cases for the

1 variance components (σ_e^2 (var_e), σ_g^2 (var_g), σ_b^2 (var_b)), genotype main effects **g**
 2 (**g**), genotype slope **b** (**b**) and the function predictor \hat{y} (yhat). There are many high
 3 peaks in the trace plot of σ_h^2 (var_h), which indicates that the distribution of σ_h^2 is
 4 skewed (this is also self-evident in the density plot). This should be expected since there
 5 are only four levels of environment effect and scaled inverse chi-square distribution with
 6 few degrees of freedom is highly skewed. Figure 3 reproduces the trace plot of the
 7 variance components (var_e, var_g, var_b, var_h).

8 The mixing for the intercept μ and the environment effects (the entries of **h**) can
 9 be slow in multiplicative models (e.g., Shariati et al., 2009). Therefore, the user is
 10 advised to check convergence to the posterior distribution and the magnitude of Monte
 11 Carlo standard errors. Convergence to the posterior distribution can be assessed
 12 graphically using a trace plot for single or more formally multiple chains. Figure 4
 13 reproduces the trace plot for intercept μ and the first two elements of environment effect,
 14 $h[1]$ and $h[2]$, in all cases we used samples from model GibbsA. From Figure 4, we
 15 can see that even the mixing of $h[1]$ and $h[2]$ is slow, when running 50,000 iterations,
 16 the chain has converged to relative constant sample means. The Time-series standard
 17 error for the sample means of $h[1]$ and $h[2]$ are both around 0.0065, which is at a
 18 reasonable level (obtained by line 6 of Box 7). An example of how to assess
 19 convergence using multiple chains is provided in the supplemental materials.

20

Box 7: Plot Gibbs samples

```
1 load("GibbsIsamps.rda")
2 plot(samps, ask=T)
3 load("GibbsAsamps.rda")
```

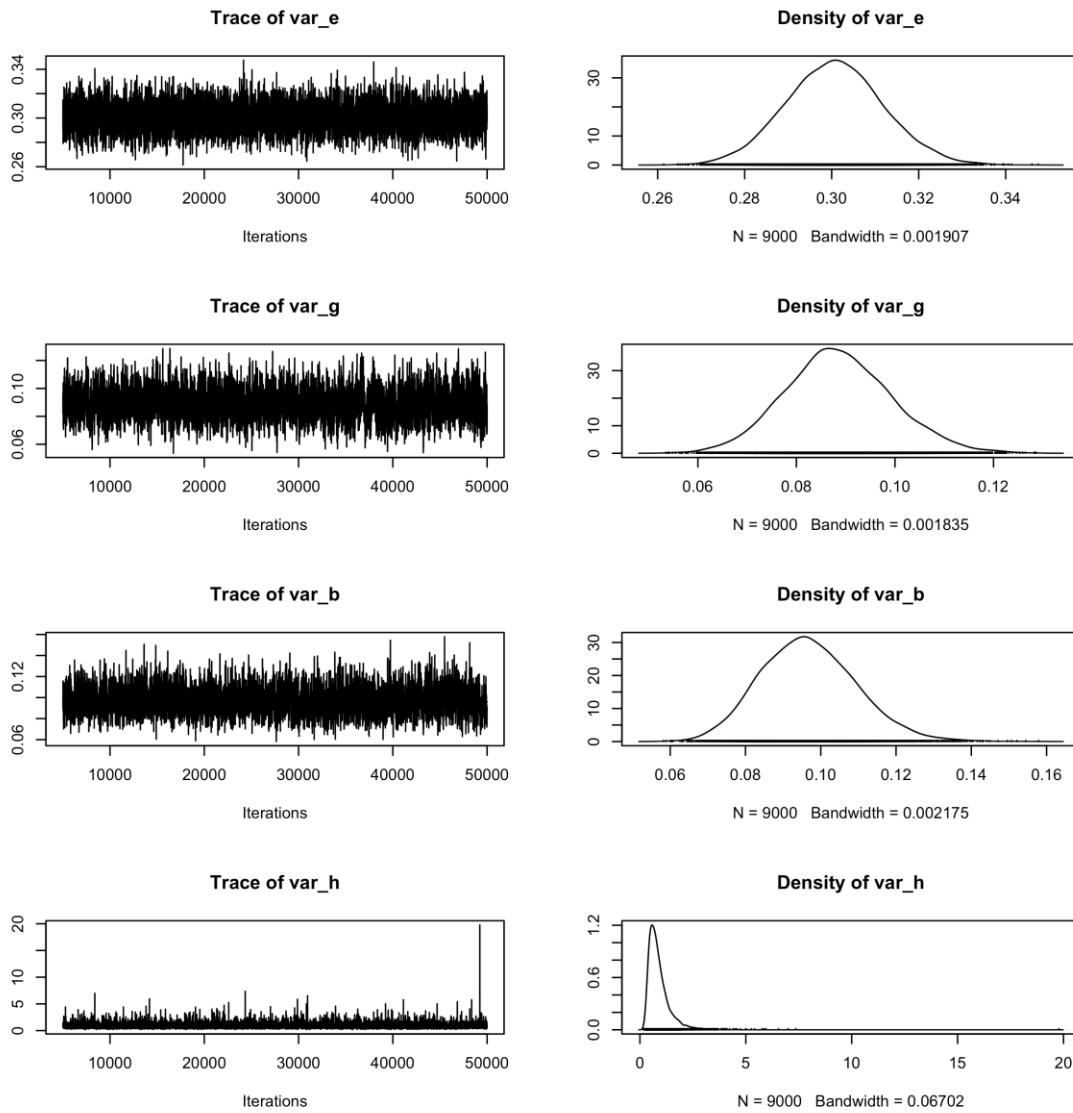
```

4 plot(samps,ask=T)
5 autocorr(samps[[1]][,"var_e"])
6 summary(samps)

```

1

2



3

4

Figure 3. Trace and density plot of variance components from GibbsA.

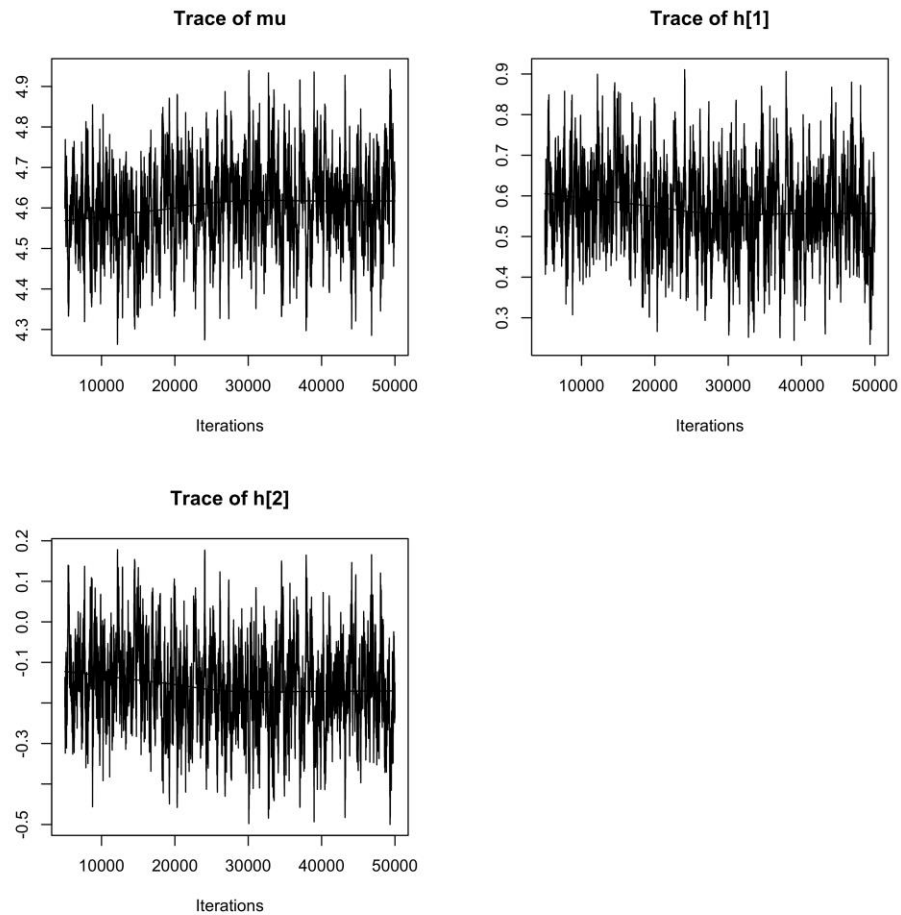


Figure 4. Trace plot of the intercept (μ) and the first two levels of environment effects ($h[1]$) and $h[2]$) from GibbsA.

Example 2: Assessment of prediction accuracy in testing data sets

Example 1 suggests that the OLS method fitted the training data better than the Bayesian models; this is expected because shrinkage reduces fitness to the data used to train a model. However, better model fitness does not necessarily imply higher prediction accuracy in validation data sets. In the following example we illustrate how to use the FW for assessment of prediction accuracy using cross-validation.

To assess the ability of different models for predicting new data, we modified the code in Box 3 by setting NA to randomly selected entries of the phenotypic vector (i.e.,

1 one record out of four per line was randomly selected and labeled as NA; see, code in
2 lines 1-5 in Box 8). The FW package produces estimates and predictions for all the lines,
3 environments and all the entries of the phenotypic vector, those that had observed values
4 and those that had NA. Therefore, predictions for entries with masked phenotypes can be
5 used to assess prediction accuracy in validation data sets (see lines 17-19 in Box 8). We
6 repeated the code in Box 8 for 100 times and generated 100 random partitions of the data
7 into training and testing sets. Each partition renders an estimate of prediction accuracy for
8 each of the models.

Box 8: correlation between y and \hat{y} for training and validation data sets.

```

1 yNA=y
2 seed=12345; set.seed(seed)
3 #randomly masking one environment for each variety
4 whichNa=seq(from=0,to=2392,by=4)+sample(1:4,size=599,replace=T)
5 yNA[whichNa]=NA
6
7 OLS=FW(y=yNA,VAR=VAR,ENV=ENV, method="OLS")
8 GibbsI=FW(y=yNA,VAR=VAR,ENV=ENV,
9 method="Gibbs",seed=seed,nIter=50000, burnIn=5000)
10 GibbsA=FW(y=yNA,VAR=VAR,ENV=ENV,
11 method="Gibbs",A=wheat.G,seed=seed,nIter=50000,burnIn=5000)
12
13 cor(y[-whichNa],OLS$yhat[-whichNa,])
14 cor(y[-whichNa],GibbsI$yhat[-whichNa,])
15 cor(y[-whichNa],GibbsA$yhat[-whichNa,])
16
17 cor(y[whichNa],OLS$yhat[whichNa,])
18 cor(y[whichNa],GibbsI$yhat[whichNa,])
19 cor(y[whichNa],GibbsA$yhat[whichNa,])

```

9

10 The mean correlation (of the 100 replicates) between phenotypes and predictions
11 in the training data set (i.e., for the entries of y that did not have missing values) follows
12 the same patterns as in Example 1, where OLS fitted the data best: 0.95 for OLS, 0.89 for

1 GibbsI and 0.86 for GibbsA. However, The mean prediction correlation (of the 100
2 replicates) for the entries of the validation set has reversed orders: 0.61 for OLS, 0.77 for
3 GibbsI and 0.80 for GibbsA.

4 In Figure 5, we plotted the estimated prediction correlation between predictions
5 and observations in training (1st row of plots) and testing (2nd row of plots) data sets.
6 Plots in the 1st, 2nd, and 3rd column correspond to comparisons of: OLS vs. GibbsI, OLS
7 vs. GibbsA and GibbsI vs. GibbsA, respectively. Within each plot each point represents
8 the accuracy obtained in a partition for the models represented in the vertical and
9 horizontal axis. Points above (below) the 45-degree line indicate higher (lower) accuracy
10 of the model in the vertical axis, relative to the one in the horizontal axis. We observed
11 that OLS always fitted the data better than GibbsI and GibbsA in the training data sets;
12 however, GibbsI and GibbsA always outperformed OLS by a sizable margin in terms of
13 prediction accuracy in testing data sets. Finally, incorporating genetic information
14 (GibbsA) always led to higher prediction accuracy than models that assumed
15 independence between lines (GibbsI).

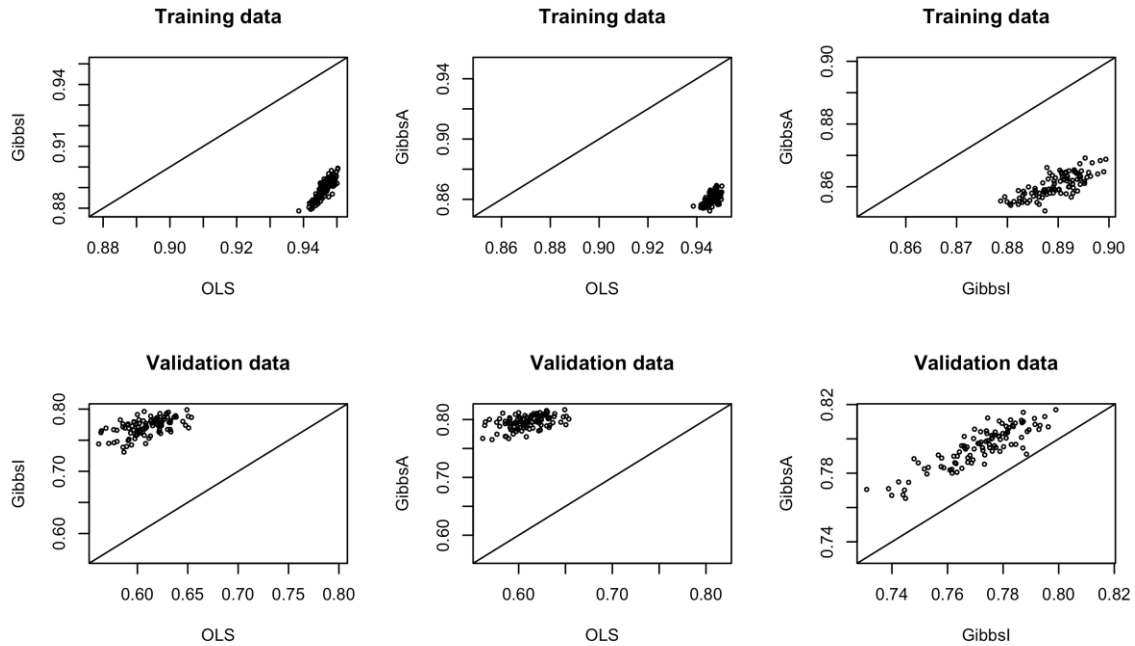


Figure 5. Prediction accuracy for training and validation sets for the three methods implemented in Box 8.

We also noted in Table 4 that the correlations (here we reported results only for the first replicate) for the parameter estimates among different models reduced compared to Example 1 due to the missing values. For example, the correlation for the estimated \mathbf{b} among different models reduced to 0.85 between OLS and GibbsI, 0.64 between OLS and GibbsA, and 0.79 between GibbsI and GibbsA.

Table 4. Pearson's product-moment correlation between parameter estimates derived by each of the three methods implemented in Box 7 (results from the first replicate only).

	OLS-GibbsI	OLS-GibbsA	GibbsI-GibbsA
\hat{h}	1.00	1.00	1.00

\hat{b}	0.85	0.64	0.79
\hat{g}	0.96	0.73	0.77
\hat{y}	0.91	0.87	0.97

1 **Computation time for 599 wheat lines**

2 We ran the FW function in an Intel Core i7 1867 MHz Processor (R was executed
3 in a single thread) with 16 GB of RAM memories. We recorded the memory and time
4 usage for Gibbs methods with 50000 iterations. With the full dataset (599 varieties, 2396
5 observations) the process used approximately 50 M of RAM memory for GibbsA, 17 M
6 of RAM for GibbsI and 153 M for OLS. The time needed to finish the process was: 11
7 minutes for GibbsA, 3 minutes for GibbsI and 2 seconds for OLS.

8 **Concluding Remarks**

9 The FW package allows fitting Finlay-Wilkinson regression with ordinary least
10 square method and Bayesian method. For Bayesian method, covariance matrix among
11 varieties and environments can be included in the model. The interface allows the user to
12 fit the models (e.g. OLS versus Gibbs) and visualize the results easily. The algorithms for
13 Gibbs Sampler are implemented in C and the speed is high. The package also provided
14 flexibility for changing the hyper-parameters and model output.

15 For incomplete/unbalanced experimental design the Bayesian approach is
16 expected to have better statistical performance and prediction accuracy than the
17 traditional two-step OLS method. Furthermore, the Bayesian models implemented in
18 FW allows incorporating pedigree, marker information as well as modeling

1 **environment covariance.** A cross-validation study based on real wheat data confirmed
2 those expectations; indeed, the Bayesian method incorporating relationships between
3 lines had a prediction accuracy that was 30% greater than the two-steps OLS method.

4 **Acknowledgements**

5 We thank the collaborators in national agricultural research institutes who carried
6 out the Elite Spring Wheat Yield Trials (ESWYT) and provided the phenotypic data
7 analyzed in this article. GDLC and LL received financial support from NIH grants
8 R01GM101219 and R01GM099992 and from Arvalis. GDLC received financial support
9 from Arvalis and CIMMYT.

1 **References**

- 2 Casella, G., and E. I. George, 1992. Explaining the Gibbs sampler. The American
3 Statistician 46: 167-174.
- 4 Copas, J. B., 1983. Regression, prediction and shrinkage. Journal of the Royal Statistical
5 Society. Series B (Methodological) 45: 311-354.
- 6 Crossa, J., G. de Los Campos, P. Pérez, D. Gianola, J. Burgueño, J. L. Araus, D.
7 Makumbi, R. P. Singh, S. Dreisigacker, and J. Yan, 2010. Prediction of genetic
8 values of quantitative traits in plant breeding using pedigree and molecular
9 markers. Genetics 186: 713-724.
- 10 Finlay, K., and G. Wilkinson, 1963. The analysis of adaptation in a plant-breeding
11 programme. Crop and Pasture Science 14: 742-754.
- 12 Frank, L. E., and J. H. Friedman, 1993. A statistical view of some chemometrics
13 regression tools. Technometrics 35: 109-135.
- 14 Geman, S., and D. Geman, 1984. Stochastic relaxation, Gibbs distributions, and the
15 bayesian restoration of images. IEEE Transactions on Pattern Analysis and
16 Machine Intelligence 6: 721-741.
- 17 Gregorius, H.-R., and G. Namkoong, 1986. Joint analysis of genotypic and
18 environmental effects. Theoretical and Applied Genetics 72: 413-422.
- 19 Jarquín, Diego, José Crossa, Xavier Lacaze, Philippe Du Cheyron, Joëlle Daucourt *et al.*,
20 2014. A reaction norm model for genomic selection using high-dimensional
21 genomic and environmental data. Theoretical and applied genetics 127: 595-607.

- Perkins, J. M., and J. Jinks, 1968. Environmental and genotype-environmental components of variability III. Multiple lines and crosses. *Heredity* 23: 339-356.
- Plummer, M., N. Best, K. Cowles, and K. Vines, 2006. CODA: convergence diagnosis and output analysis for MCMC. *R News* 6: 7-11.
- Pérez, P., G. de Los Campos, J. Crossa, and D. Gianola, 2010. Genomic-enabled prediction based on molecular markers and pedigree using the bayesian linear regression package in R. *The Plant Genome* 3: 106-116.
- R Development Core Team, 2011. R: a language and environment for statistical computing. R Foundation for Statistical Computing, Vienna and Austria.
- Shariati, M., I. Korsgaard, and D. Sorensen, 2009. Identifiability of parameters and behaviour of mcmc chains: a case study using the reaction norm model. *Journal of Animal Breeding and Genetics* 126: 92-102.
- Su, G., P. Madsen, M. S. Lund, D. Sorensen, I. R. Korsgaard, and J. Jensen, 2006. Bayesian analysis of the linear reaction norm model with unknown covariates. *Journal of Animal Science* 84: 1651-1657.
- Walsh, J. B., and M. Lynch, 2014. Chapter 44 Selection and G x E: Advanced Topics. *In* *Evolution and Selection of Quantitative Traits: II. Advanced Topics in Breeding and Evolution*. Available at http://nitro.biosci.arizona.edu/zbook/NewVolume_2/pdf/Chapter44.pdf (verified 05 Dec. 2015).