

Bayesian Ridge Regression Shows the Best Fit for Ssr Markers in Psidium Guajava Among Bayesian Models

Flávia Alves Silva (✉ flavia_uems@hotmail.com)

Universidade Estadual do Norte Fluminense Darcy Ribeiro (UENF)

Alexandre Pio Viana

Universidade Estadual do Norte Fluminense Darcy Ribeiro (UENF)

Caio Cezar Guedes Corrêa

Universidade Estadual do Norte Fluminense Darcy Ribeiro (UENF)

Eileen Azevedo Santos

Universidade Estadual do Norte Fluminense Darcy Ribeiro (UENF)

Julie Anne Vieira Salgado Oliveira

Universidade Estadual do Norte Fluminense Darcy Ribeiro (UENF)

José Daniel Gomes Andrade

Universidade Estadual do Norte Fluminense Darcy Ribeiro (UENF)

Rodrigo Moreira Ribeiro

Universidade Estadual do Norte Fluminense Darcy Ribeiro (UENF)

Leonardo Siqueira Glória

Universidade Estadual do Norte Fluminense Darcy Ribeiro (UENF)

Research Article

Keywords: SAM, Bayesian alphabet, cross-validation, genetic correlation, heritability

Posted Date: March 13th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-284606/v1>

License: © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Abstract

Markers are an important tool in plant breeding, which can improve conventional phenotypic breeding, generating more accurate information outcoming better decision making. This study aimed to apply and compare the fit of different Bayesian models BRR, BayesA, BayesB, BayesB (setting the value from very low to $\pi = 10^{-5}$) and BayesC and Bayesian Lasso (LASSO) for predictions of the genomic genetic values of productivity and quality traits of a guava population. A randomized block design with two replications was used. Seventeen full-sib families were evaluated. Fruit mass, pulp mass, soluble solids content, number of fruits, and production per plant were measured. These variables were used in the genomic prediction with SSR markers, obtained through the CTAB extraction method with 200 primers. The Bayesian ridge regression model showed the best results for all variables and was chosen to predict the individuals' genomic values according to the cross-validation data. A good stabilization of the Markov and Monte Carlo chains was observed with the mean values corresponding to the observed phenotypic means. Heritabilities showed good predictive accuracy. The model showed strong correlations between some variables, allowing indirect selection.

Keywords: SAM; Bayesian alphabet, cross-validation, genetic correlation, heritability.

Introduction

Tropical fruits have a great commercial value worldwide because, besides being widely consumed in the countries that produce them, they are highly appreciated and with a great added value around the world ¹. One of these perennial fruits is the guava tree (*Psidium guajava* L.), which is gaining space on the market due to the increasingly efficient selection methods for improving the species.

One of these methods is the selection of superior individuals with the help of molecular markers, such as genomic selection. This method characterizes the ideal association between conventional breeding based on phenotypic observations and modern molecular techniques currently available. Its use has a great impact on breeding programs ². This is because its main objective is to obtain more accurate and precise estimates, which allows for better planning.

1 However, the breeder has available several statistical models to associate the marks with the
2 phenotypes, which makes it a challenge to choose a suitable model for the response of the species and marks.
3 Recently, among these models, Bayesian approaches have gained a lot of prominence with the advent of
4 computational power. With a Bayesian approach, the effects of the markers can be estimated together to predict
5 the genomic values for a quantitative trait without making the previous selection in the panel of markers ³. This
6 Bayesian genomic selection has as main advantages the inclusion of a priori information in the model, besides
7 generating more accurate credibility intervals ⁴.

8 Accuracy varies between models of genomic selection, according to their assumptions and treatments
9 of the effects of the markers. For example, it was identified that Bayesian models (*Bayesian LASSO - BL*) and
10 ridge regression models (BRR) showed superior performance for traits controlled by additive genetic effects ⁵.

11 Among the available Bayesian approaches for estimating genomic values in plant breeding, we can
12 mention LASSO Bayesian - BL that combines both selection and variable contraction methods. Advantageous
13 concerning the most common method that does not use variable selection. It has an exponential priori in the
14 variance of the markers, resulting in a double exponential distribution. The double exponential distribution has
15 a high mass density at zero, and heavier priori tails compared to a Gaussian distribution ^{6,7}. Bayesian ridge
16 regression -BRR induces homogeneous shrinkage of all marker effects to zero and produces a Gaussian
17 distribution of marker effects ⁸.

18 BayesA uses an inverse-chi-square (x^2) in the variance of the markers, producing a scaled t distribution
19 for the effects of the markers. Similar to BL and unlike BRR, it shrinks the markers with small effects to values
20 close to zero, and the markers with greater effects are maintained. The final distribution of the marks shows a
21 higher peak of mass density close to zero compared to the double exponential distribution ^{6,9}. BayesB is similar
22 and uses an inverse x^2 but uses shrinkage and selection methods of the variable. And when the priori parameter
23 $\pi = 0$, it is like BayesA ¹⁰. BayesC also applies the shrinkage and selection methods of variable and generates a
24 Gaussian distribution of the effects of the markers. BayesB and BayesC consist of close to zero density in the
25 distribution when using low priori ^{6,11,12}.

26 For the breeder, finding out which model best fits his object of study is of paramount importance for
27 the planning of the breeding program. For guava, there is not yet a study looking for which model is best applied

1 to the association of marks, although primers for simple-sequence repeats (SSR) have also been applied, as
2 observed in Dinesh, et al. ¹³.

3 This study aimed to apply and compare the fit of different Bayesian models BRR, BayesA, BayesB,
4 BayesB (setting the value from very low to $\pi = 10^{-5}$) and BayesC and Bayesian Lasso (LASSO) for predictions
5 of the genomic genetic values of productivity and quality traits of a guava population.

6 **Material and methods**

7 **Genetic material**

8 The data used in this study were obtained in the experiment carried out in the Guava Breeding Program
9 of UENF and are in accordance with the institutional guidelines for carrying out the experiments. The
10 experimental area was located at the Antônio Sarlo Technical and Agricultural School, in Campos dos
11 Goytacazes, Rio de Janeiro, Brazil. In the experimental field, a complete block design with two replications per
12 plot was used. Each plot contained one of the seventeen guava segregating families with twelve plants.

13 The segregating families were obtained by crossings between the accessions, that were established
14 considering information on genetic diversity obtained by Pessanha, et al. ¹⁴. Twelve families were selected and
15 from them we selected the best individuals from each family based on the work of Silva, et al. ¹⁵ to apply the
16 markers were:

17 UENF 1834 × UENF 1833 (12 individuals);

18 UENF 1831 × UENF 1830 (12 individuals);

19 UENF 1831 × UENF 1832 (1 individual);

20 UENF 1833 × UENF 1832 (11 individuals);

21 UENF 1834 × UENF 1839 (1 individual);

22 UENF 1835 × UENF 1834 (16 individuals);

23 UENF 1836 × UENF 1835 (15 individuals);

24 UENF 1833 × UENF 1836 (2 individuals);

25 UENF 1831 × UENF 1835 (10 individuals);

26 UENF 1833 × UENF 1835 (5 individuals);

27 UENF 1834 × UENF 1837 (5 individuals);

28 UENF 1832 × UENF 1835 (6 individuals).

1 Five explanatory variables were measured for each individual: fruit mass (FM), pulp mass (PM),
2 soluble solids content (SSC), number of fruits per plant (NF), and production per plant (PROD). Five
3 observations of all variables were obtained, except for NF and PROD, for which one observation was carried
4 out per individual.

5 **DNA extraction and quantification**

6 DNA extraction was carried out from young leaves collected individually from each plant, using the
7 standard CTAB method with modifications ¹⁶. Then, the DNA was quantified by analysis on 1% agarose gel on
8 TAE 1X buffer (Tris, Sodium Acetate, EDTA, pH 8.0), using the Lambda marker (λ) of 100 bp ($100 \text{ ng} \cdot \mu\text{L}^{-1}$)
9 (Invitrogen, USA), by comparing the bands. For this procedure, the samples were stained using the mixture of
10 Gel, RedTM, and Blue Juice (1:1), and the image was captured by the MiniBis Pro photodocumentation system
11 (Bio-Imaging Systems). Subsequently, the DNA samples were diluted to a working concentration of $10 \text{ ng} \cdot \mu\text{L}^{-1}$.
12

13 **Primer Screening**

14 Two-hundred SSR markers were tested ¹⁷ in five guava individuals to identify polymorphic loci. After
15 screening, a set of 44 polymorphic primers was selected for the amplification reactions in a population of 96
16 individuals from the field experimente with phenotypic data. These individuals were selected for their
17 performance after seven years of phenotypic data, and represent the individuals who will proceed to the next
18 stages of the breeding program.

19 **Polymerase chain reaction (PCR)**

20 The PCR reactions were carried out in thermocyclers from Applied Biosystems/Veriti 96 well, in a 38
21 cycle program, obeying the following temperatures and time: 94 °C for one minute (initial denaturation), 94 °C
22 for two minutes (cyclic denaturation), the specific temperature of each initiator, in °C, for one minute
23 (annealing), 72 °C for three minutes (cyclic extension), 72 °C for 10 minutes (final extension), and 4 °C. The
24 final volume was 13 μL of each sample, being: 2 μL of DNA ($10 \text{ ng}/\mu\text{L}$), 1.50 μL of 10X Buffer (NH_4SO_4), 1.5
25 μL of MgCl_2 (25 mM), 1.5 μL of dNTPs (2 mM), 1 μL of primer (R+F) (5 μM) and 0.12 μL of Taq-DNA
26 polymerase (5 U/ μL) (Invitrogen, Carlsbad, Califórnia, EUA). The amplification products were separated on
27 4% Metaphor agarose gel, stained with GelRedTM, and visualized through the MiniBis Pro photo-
28 documentation system (Bio-Imaging Systems).

1 Statistical analysis

2 The phenotypes of each individual were used as a response variable in the genomic predictions, using
3 the following models: Bayesian Ridge regression (BRR – Bayesian Ridge regression), BayesA, BayesB,
4 BayesB (setting the very low value of π , 10^{-5}), BayesC and Bayesian Lasso (Bayesian Lasso - BL, assuming
5 the marginal distribution as double exponential prior to the effects of markers). The general model for genomic
6 predictions can be described in the matrix form as:

$$7 \quad y = \mu + Wg + e \quad (1)$$

8 where: y is the vector of the observations for each characteristic, μ is the general mean, g is the vector
9 with the effects of the markers, whose assumptions depend on the model used, W is the matrix of the genotypes
10 (coded as 0, 1, and 2) of each plant for each marker and e is the vector of the residues.

11 A complete description of the calculation of heritability and the specifications of the probability
12 distributions of the general model effects above, for the use of Bayesian methods, can be found in Pérez and de
13 Los Campos ¹⁸. All Bayesian analyzes were performed in the BGLR package ¹⁸ of the R software ¹⁹, with the
14 BGLR function adjusted for 1E6 iterations with the first 2E5 cycles discarded as burn-in and thin assuming the
15 value 4.

16 The models were compared based on the Deviance Information Criterion (DIC) proposed by
17 Spiegelhalter, et al. ²⁰. The DIC can be described as follows $DIC = D(\hat{\theta}) + 2p_D$, in which the first term is a
18 Bayesian model adjustment measure ($D(\hat{\theta})$), which is defined as the a posteriori mean of deviance and the
19 second component (p_D) measures the complexity of the model through the effective number of parameters.
20 Posterior probabilities of the models were calculated using the approximation presented by Wilberg and Bence
21 ²¹ to facilitate the interpretation of DIC values in terms of the superiority of one model over the other, in which
22 it is given by:

$$23 \quad p(M_t \vee l) = \frac{\exp(-\Delta_t/2)}{\sum_{t=1}^6 \exp(-\Delta_t/2)}, t = 1,2,3,4,5,6 \quad (1)$$

24 where: $p(M_t \vee l)$ is the a posteriori probability of model t , Δ_t is the difference between the DIC of
25 model t and the model with the lowest DIC.

26 Results

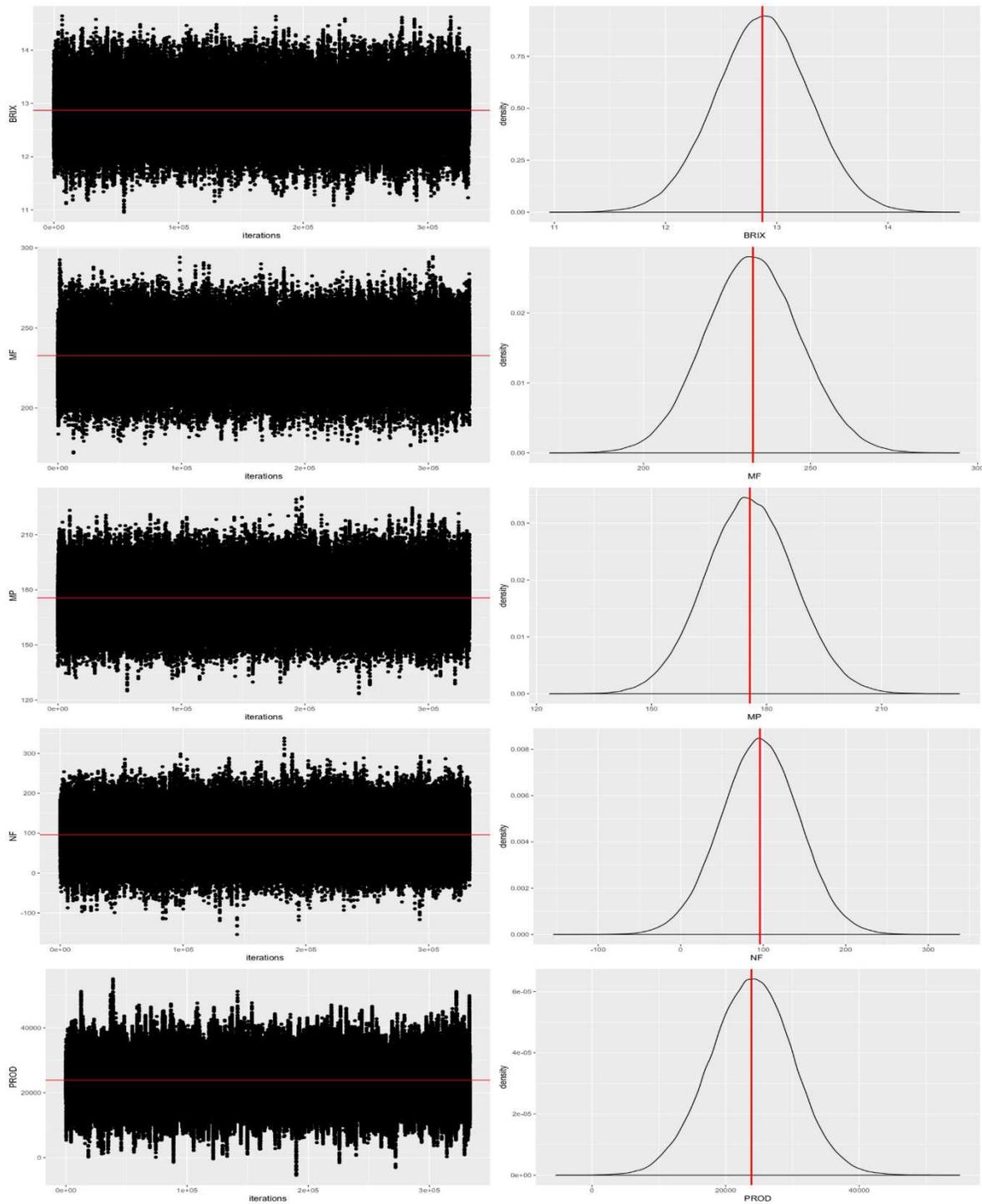
27 Six Bayesian models were applied to detect the effect of the markers along with phenotypic data from
28 a guava population. In the modeling process, cross-validation with eight folds was used to obtain some adjust

1 parameters of the models in all folds (Table 1). Among the models used, the Bayesian Ridge Regression model
2 - BRR presented the lowest mean value considering a comparative adjustment value ($<DIC - Deviance$
3 *Information Criterion*) according to the parameters used in the variable soluble solids content ($^{\circ}$ BRIX).

4 The DIC is particularly useful in problems of Bayesian selection models, where the posterior
5 distributions of the models were obtained by the Markov Chain Monte Carlo simulation (MCMC). DIC is an
6 asymptotic approach as the sample size becomes large, like the AIC. It is only valid when the posterior
7 distribution is approximately normal multivariate. Thus, the chain convergences and the posterior distribution
8 (normal distribution) were verified for all variables in the BRR model (Fig 1).

9 A good stabilization of the Markov and Monte Carlo chains was observed with the mean values
10 corresponding to observed phenotypic means. The posterior density curves of the chains showed normal
11 distribution in all variables. Therefore, it is possible to use DIC values to select the models safely.

12 Deviations (Δ) of information criteria were also obtained for each variable concerning the lowest value,
13 assumed as the model that presented the best fit to the data. From these parameters, auxiliaries were also
14 obtained in the classification of models as values of posterior adjustment probability of the model (W_{prob}) and
15 the evidence ratio (ER) for the models. All adjust parameters of the BRR model were superior to the other
16 Bayesian models used for the SSC variable.



1
 2 **Fig 1.** Markov and Monte Carlo chains with mean values (red line) and distribution curve for five variables
 3 observed in guava, generated to relate SSR marks to phenotypic observations.

4 Besides the adjustment values, for the model choice, we consider the model's ability to predict the
 5 phenotypic values of a sub-sample with random individuals, in each fold of the cross-validation. The mean

1 values of the predicted correlation and the observed phenotypes (r), together with a probability value of r , had
 2 no linear correlation (Table 1).

3 For the SSC variable, the BRR model also showed the highest r value with the lowest probability,
 4 being a consistent correlation between the subsamples. The other models performed very similarly, except for
 5 the BayesL model, where a discrepant DIC value was observed, and the BayesB2 model, which despite showing
 6 a good fit with a similar DIC, presented a low predictive capacity with $r = 0.35$ concerning BRR with $r = 0.65$.

7 A similar result in the model's adjustment and prediction criteria was observed for the other variables,
 8 such as the number of fruits per plant (NF). In NF, adjusted values of the very similar models close to 980.19
 9 (DIC) were observed, with the BRR model chosen by the best predictive capacity with $r = 0.82$ for 0.65 for
 10 FM, 0.64 for PM, and 0.84 for PROD.

11 With the model adjustment criteria very close between the models used and great differences between
 12 the predictive power of each model within the variables, it was possible to observe that choosing a value of π
 13 for the BayesB model caused an overfit of the model. It was observed that the predictive power of the BayesB
 14 model, in most cases, presented the worst results (r).

15 **Table 1.** Adjustment quality of six Bayesian models: BL, BRR, BayesA, BayesC, BayesB, and BayesB with π
 16 = $1e-5$ (BayesB2) to associate SSR markers and phenotypic data in *P. guajava* in the variables of soluble solids
 17 content, fruit mass, pulp mass, number of fruits per plant and production per plant. The bias values were
 18 obtained by eight-fold cross-validation (88% of the data for training and 12% for validation), in the same sample
 19 sets for each model.

	DIC	Δ	Wprob	ER	r	p-value
Soluble solids content						
BRR	1348.95	0.00E+00	5.68E-01	1.00E+00	0.65	1.76E-11
BayesA	1352.36	3.41E+00	1.03E-01	5.50E+00	0.65	1.82E-11
BayesL	1455.71	1.07E+02	3.73E-24	1.52E+23	0.65	3.34E-11
BayesC	1352.66	3.71E+00	8.90E-02	6.38E+00	0.65	2.50E-11
BayesB	1353.17	4.22E+00	6.87E-02	8.26E+00	0.65	3.29E-11
BayesB2	1351.35	2.40E+00	1.71E-01	3.31E+00	0.35	3.00E-10

Fruit mass

BRR	4330.52	7.03E-01	2.58E-01	1.42E+00	0.65	5.53E-12
BayesA	4332.11	2.29E+00	1.16E-01	3.15E+00	0.65	8.88E-12
BayesL	4377.44	4.76E+01	1.67E-11	2.20E+10	0.64	1.20E-11
BayesC	4331.05	1.23E+00	1.98E-01	1.85E+00	0.65	7.59E-12
BayesB	4333.38	3.57E+00	6.16E-02	5.95E+00	0.65	1.04E-11
BayesB2	4329.82	0.00E+00	3.66E-01	1.00E+00	0.52	2.46E-12

Pulp mass

BRR	4179.22	8.37E-01	2.64E-01	1.52E+00	0.64	3.04E-11
BayesA	4181.10	2.71E+00	1.03E-01	3.89E+00	0.64	5.58E-11
BayesL	4224.53	4.61E+01	3.83E-11	1.05E+10	0.62	7.56E-11
BayesC	4180.16	1.78E+00	1.65E-01	2.43E+00	0.63	4.36E-11
BayesB	4181.99	3.61E+00	6.60E-02	6.08E+00	0.63	6.04E-11
BayesB2	4178.39	0.00E+00	4.01E-01	1.00E+00	0.50	1.28E-11

Number of fruits

BRR	980.19	6.57E-01	1.71E-01	1.39E+00	0.82	4.36E-13
BayesA	981.18	9.87E-01	1.45E-01	1.64E+00	0.79	8.92E-13
BayesL	980.75	5.58E-01	1.79E-01	1.32E+00	0.76	4.31E-13
BayesC	981.36	1.17E+00	1.32E-01	1.79E+00	0.80	2.88E-13
BayesB	980.15	0.00E+00	2.37E-01	1.00E+00	0.79	2.39E-13
BayesB2	981.31	1.12E+00	1.35E-01	1.75E+00	0.73	8.51E-13

Production per plant

BRR	1825.73	8.71E-01	1.99E-01	1.55E+00	0.84	2.02E-13
BayesA	1826.51	1.64E+00	1.36E-01	2.27E+00	0.82	5.34E-13
BayesL	1828.81	3.95E+00	4.28E-02	7.20E+00	0.77	4.08E-13
BayesC	1825.92	1.06E+00	1.81E-01	1.70E+00	0.82	1.47E-13
BayesB	1826.55	1.69E+00	1.33E-01	2.33E+00	0.81	3.53E-13
BayesB2	1824.86	0.00E+00	3.08E-01	1.00E+00	0.74	1.10E-12

1 DIC = deviance information criterion; Del (delta) = difference between the highest and the lowest DIC value;
 2 Wprob = posterior probability model; ER = evidence ratio; Error = error attributed to Wprob; r = correlation
 3 between predicted by the model and reserved validation data; p-value = significance of the correlation.

4 It is worth mentioning that the Bayesian models take considerable time to be executed. Even with the
 5 advancement of computational power, the resolution of more complex models requires a longer processing
 6 period. This is widely known information, but little measured, which must be considered when choosing the
 7 model. In this study, the time invested in solving the Bayesian models was measured by repeating each chain
 8 ten times in a loop (Table 2).

9 **Table 2.** Estimates of time averages for solving different Bayesian models with 10^6 iterations, burn-in of 10^4 ,
 10 and chain sampling equal to 3 The calculations were performed in the R software (version 3.5.1) with the BGLR
 11 package (version 1.0.8).

Model	Processing time*
Bayesian Ridge Regression - BRR	15.3 hours (+- 0.15)
BayesA	15.35 hours (+- 0.10)
BayesB	15.32 hours (+- 0.13)
BayesB $\pi = 10^{-5}$	15.3 hours (+- 0.10)
BayesC	15.5 hours (+- 0.2)
Lasso Bayesian	14.95 hours (+- 0.5)

12 * A 2.7 GHz Intel I7-7500U processor core was used.

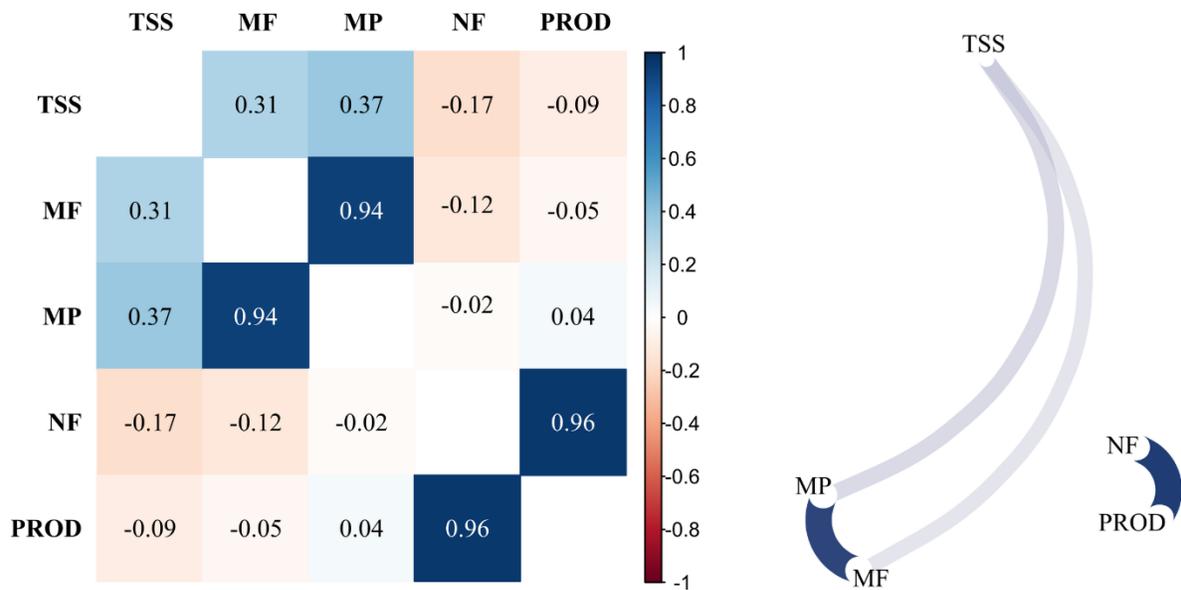
13 Narrow-sense heritability values were estimated for the variables observed in guava with the model
 14 that showed the best performance (Table 3). The extremes of the values were for the soluble solids content (TSS
 15 = 0.32) with the highest observed heritability value, and the number of fruits per plant (NF = 0.07) had the
 16 lowest heritability value. In general, heritability values were low but accompanied by deviation and accuracy
 17 measures; they can provide more accurate estimates for the advancement of generations in the breeding
 18 program. The values of the heritability deviation measure were low. This indicates more precise values for
 19 heritability, as opposed to estimates of heritability obtained punctually, as is commonly done. Error estimates
 20 of heritability were obtained with the estimate of heritability in each iteration of cross-validation, thus being
 21 estimated in several subsets that represent the population.

In predictive accuracy, high values were observed, which is a good indication that the estimated heritability represents the population very well. In particular, the predictive accuracy value for the PM variable, estimated at 0.9708. However, when observing the predictive accuracy of heritability of PROD, the variable of main interest, a value close to 0.51, was obtained, which is low. Very similar results were observed for fruit mass and pulp mass. The heritability values were 0.1581 and 0.1478 from FM and PM, respectively. The standard deviation of heritability was also close and low. Only the predictive accuracy was better in PM than FM, indicating that the volume inside the fruits depends less on the size of the fruit, being more random or influenced by another factor not observed in this experiment.

Table 3. Predictive accuracy and standard deviation of heritability for soluble solids content (TSS), fruit mass (FM), pulp mass (PM), number of fruits per plant (NF) and production per plant (PROD) observed in guava (*Psidium guajava*), estimated using a model with SSR markers and Bayesian ridge regression - BRR.

	H2	Standard Deviation	Predict Accuracy
SSC	0.3261	0.0706	0.6302
FM	0.1581	0.0215	0.8779
PM	0.1478	0.0310	0.9708
NF	0.0732	0.0146	0.6939
PROD	0.1058	0.0154	0.5095

With the matrix of the individuals' marks and the weight that each marker received in the Bayesian ridge regression model, the individuals' genetic values of the traits were estimated, and the genetic correlation matrix between the traits was obtained (Fig 2). A high linear correlation was observed between PM and FM (0.9393) and between NF and PROD (0.9641). A correlation was also observed between soluble solids content and two variables of the fruit, with a value of 0.3060 between SSC and FM, and 0.3705 between SSC and PM. It was also observed that SSC showed a negative correlation with PROD and NF, but there were no significant correlations.



1

2 **Fig 2.** Genetic correlation between the soluble solids content (TSS), fruit mass (FM), pulp mass (PM), number
 3 of fruits per plant (NF), and production per plant (PROD) observed in guava (*Psidium guajava*), estimated using
 4 a model with SSR markers and Bayesian ridge regression - BRR.

5 **Table 4.** List of selected individuals who presented positive genetic values in the variables soluble solids content
 6 (TSS), fruit mass (FM), pulp mass (PM), number of fruits per plant (NF), and production per plant (PROD).

Individual	TSS	FM	PM	NF	PROD
B1F15P12	0.76	10.68	8.45	24.57	4507.96
B2F8P4	0.19	6.54	4.12	13.11	3115.68
B2F1P4	0.64	4.99	12.08	14.53	2958.55
B1F15P10	2.51	18.14	16.13	1.25	2848.35
B1F2P8	1.21	33.01	29.67	13.28	2268.82
B2F2P10	0.86	12.18	12.33	16.93	2206.82
B2F3P11	0.89	19.74	19.79	9.29	1421.55
B1F8P1	1.04	18.97	15.93	7.67	1264.91
B2F12P9	0.96	6.90	5.19	9.00	1245.41
B2F17P2	0.27	7.61	4.86	6.15	622.86

7

The individuals were classified in descending order considering the PROD.

1 The individuals contained in table 4 were selected because they present positive values in all traits.
2 However, it is possible to use a selection index if the objective is to select new individuals to compose a new
3 population within a breeding program. Families 8, 10, and 17 were the families that contained more individuals
4 in ordering the genetic values considering the production per plant. Thus, these families have low variability
5 among themselves, but with a high productive capacity, being recommended for selection and continuity in
6 trials of Value for Cultivation and Use.

7 **Discussion**

8 From a Bayesian approach, the effects of markers can be estimated together to predict the genomic
9 values for a quantitative trait without performing the marker selection. This approach is called genomic
10 selection. Several penalized and of estimation methods of Bayesian contraction are available, for example,
11 Bayesian counterparts of Ridge Regression (Ridge Regression - RR)²², Least Absolute Shrinkage and Selection
12 Operator (Least Absolute Shrinkage and Selection Operator – LASSO)²³, as well as models such as BayesA
13 and BayesB and their extensions⁹. These models are frequently tested for different crops of interest; however,
14 for guava, this information is still scarce. In this study, the performance of six Bayesian models for adjusting
15 SSR markers in guava is discussed and estimated parameters of interest to the breeder in a breeding program.

16 Although there are differences between the methods, in a priori assumptions about the effects of the
17 markers, it was observed that adjustment parameters of the models were similar. No evident difference was
18 detected for any of the traits, mainly for DIC. Thus, the models were chosen, considering not only the adjustment
19 parameters but also their predictive capacity and how they behave concerning the markers to generate the
20 regressive model.

21 BayesL produces a stronger shrinkage of regression coefficients close to zero and less shrinkage for
22 those with large absolute values, leading to a scarcer model. By other hand, BRR reduces strongly regression
23 coefficients that have large absolute values.²⁴ Thus, it was observed that BayesL presented a median
24 performance, possibly because the number of significant marks, with great effects on the model was too scarce
25 to explain the quantitative traits evaluated. Intuitively the reverse occurred with the BRR model, which
26 considered the effects of marks more, generating a model with more marks to explain traits controlled by several
27 genes. This means that the distribution of the marks was, on average, slightly less than peaks for the effects
28 research grid in the BRR model.

1 Studies that seek the best models for different species are important to direct breeding programs. For
2 example, for another perennial plant (*Passiflora edulis*), it was observed that the BayesC model was the best
3 model for several variables evaluated in this species ²⁵. This model assumes a common variance for all effects
4 of markers but also assumes that some markers do not affect ¹¹. Thus, genes with the same allelic frequency
5 probably explain the same portion of genetic variation, suggesting that several genes with few effects control
6 the variables, as the quantitative variables observed in this study. In the results, it was possible to observe that
7 this model also presented a satisfactory performance for variables in guava, being able to be chosen as an
8 alternative model.

9 Similar results between Bayesian methods such as BayesA and BayesB and other derivatives of these
10 were also observed ¹¹, as obtained in this work. This similar result was already expected since the models have
11 few variations between them. For example, BayesB and BayesA are more tolerant of the assumption of common
12 variance between the effects of the markers. A priori assumed in these models for the effect of a j^{th} marker is a
13 joint distribution with a probability π for the beta for the mark equal to zero.

14 When the BayesB model was proposed, π was suggested with a value close to 0.95 ⁹. However, with a
15 few marks, it is possible to choose lower values for π , where BayesB with π reduced to zero is equivalent to the
16 BayesA model. As possible, overfitting of the BayesB model was observed when we used a value of $\pi = 1e-5$;
17 it was forced that the marks had a high probability of influencing the variable of interest. Thus, a model was
18 obtained in which the betas referring to the brands fitted very well to explain the sub-sample in each fold of the
19 cross-validation, but failed to predict the validation sample as observed for most traits (Table 1).

20 If only the model's adjustment parameters such as the DIC, which are widely used, had been used,
21 perhaps it would not be observed that the predictive power of the BayesB model had the worst performance.
22 This highlights the importance of cross-validation. Cross-validation was used to assess how the results of one
23 statistical model resemble another set of data. For example, how an adjusted model will predict data that was
24 not used to adjust the model. Predicting the performance of genotypes with phenotypes yet to be observed (for
25 example, newly developed lines or lines that have been evaluated in a few environments) is essential in plant
26 breeding. Therefore, cross-validation appears to be a natural way to assess model performance from the
27 breeder's perspective ²⁶.

1 Simulation studies have shown that genomic selection using markers alone can adjust the model to an
2 accuracy of up to 85% ⁹. The accuracy of 85% is the correlation between the true genetic values and the
3 predicted values of individuals in the next generation. True genetic values are known only in simulation studies.
4 In the analysis of real data, the predictability of a model must be extracted from a cross-validation study. The
5 predictability obtained from cross-validation and the quality of the model's fit do not necessarily agree with
6 each other. Starting with a small number of markers, both can increase as the number of markers increases.
7 Further increasing the number of markers may continue to increase the quality of the model's fit, but
8 predictability may drop ²⁷.

9 The heritability coefficient influences the prediction of genomic genetic values, predictive capacity,
10 and association analysis across the genome. With greater heritability of phenotype, there are improvements in
11 the identification of individuals to be used as parents in the next generations, also favoring the identification of
12 regions associated with a characteristic of interest ²⁵.

13 The heritability of the TSS characteristic was the highest observed, and the value corroborates within
14 the range with a study that evaluated a large population of guava trees in India ¹³. The authors also detect a
15 correlation between this variable and fruit mass, allowing an indirect selection. It is also suggested that there
16 may be a possible non-additive effect on the genes controlling this characteristic, as they observed a phenotypic
17 variance greater than the genotypic variance. Our model showed low predictive accuracy for heritability despite
18 the higher value. Also, approximately 40% of the subsample of validation, the model presented a biased
19 prediction, also corroborating the idea of gene action with non-additive effects from this characteristic.

20 For the other traits, the heritability values were low, as expected. The values generally reported for
21 traits such as fruit mass, pulp mass, number of fruits, and production are generally close to 0.60 ²⁸. Our estimates
22 are possibly lower because they are estimates from a model that considers the effect of marks, and the usual
23 estimates are obtained from phenotypic data that have many more sources of variation, often not considered.
24 Despite being low, heritability showed good predictive accuracy in cross-validation, reaching 0.97 for the pulp
25 mass.

26 Pulp mass is strongly correlated with the fruit mass, which from the point of view of plant physiology
27 was already expected. FM and PM are variables obtained in similar ways, where one measures the mass of the
28 whole fruit, and in the other, the placenta containing the seeds is removed, a part that does not matter in the

1 processing of the fruit. Both variables showed similar heritability values of 0.14 and 0.15, which were superior
2 to the variables of interest regarding production (NF and PROD). Generally, collinearity is observed between
3 these two variables, and this collinearity is particularly interesting for studies of correlations between variables
4 in guava, which may involve modeling structural equations such as path analysis, which seek variables that can
5 be selected indirectly.

6 As the heritability is very similar in the two variables, and the genetic correlation between them is also
7 high, a program can direct the selection of individuals with higher pulp mass with the indirectly selecting based
8 on fruit mass. In the selection stages, there is a big difference in time and resource spent between just measuring
9 the mass of a set of fruits versus opening a fruit removing the placenta, and measuring the mass of the pulp.

10 In the variables number of fruits and production, heritability was very low, together with estimates of
11 predictive accuracy. Since these are also quantitative variables, usually controlled by many genes, low
12 heritability was expected. However, despite predicting the validation subsample with more than 50%, probably
13 our model was not able to capture all the effects for these traits since the model has an unsatisfactory
14 performance.

15 It is worth mentioning that these traits evaluated are highly influenced by the environment, and
16 especially by management ²⁹. For example, a common crop handling in guava trees is the pruning and a
17 subsequent thinning of new shoots that arise after pruning. This serves to control both the plant height to
18 facilitate harvesting and the number and size of fruits. Thus, the inflorescences that originated the fruits appear
19 in buds in the axils of the new shoots. If many shoots are maintained after pruning, the number of fruits tends
20 to increase, but the fruit mass is less due to the greater distribution of the available resources of the parent plant.
21 This leads us to look for a correlation between the number of fruits and production with, for example, the mass
22 of the fruit, which was not found here, or at least it is a non-linear correlation since the correlations between NF
23 and FM are close to zero (Fig 2).

24 Different genetic values were observed among the selected individuals; a possible explanation for this
25 fact is that the population has high genetic variability. This implies in the differences between the genetic values
26 of individuals, making them more pronounced, making it easier for the methods to classify individuals with
27 greater accuracy.

28 **Conclusion**

1 The Bayesian ridge regression model showed the best results and was chosen to predict the genetic
2 values of individuals in the variables soluble solids, fruit mass, pulp mass, number of fruits, and production per
3 plant. Heritability values showed good predictive accuracy. Genetic correlations were obtained to verify the
4 relationship between variables, and the model showed strong correlations between some variables, allowing the
5 indirect selection.

7 References

- 8 1 FAO, F. *Food and Agriculture Organization of the United Nations*,
9 <<http://www.fao.org/faostat/en/#data>> (2020).
- 10 2 Leon, N. d., Jannink, J. L., Edwards, J. W. & Kaeppler, S. M. Introduction to a special issue on
11 genotype by environment interaction. *Crop Science* 56, 2081-2089,
12 doi:doi.org/10.2135/cropsci2016.07.0002in (2016).
- 13 3 Mutshinda, C. M. & Sillanpää, M. J. Extended Bayesian LASSO for multiple quantitative trait loci
14 mapping and unobserved phenotype prediction. *Genetics* 186, 1067-1075,
15 doi:doi:10.1534/genetics.110.119586 (2010).
- 16 4 Xavier, A. Efficient estimation of marker effects in plant breeding. *G3: Genes, Genomes, Genetics* 9,
17 3855-3866, doi:doi.org/10.1534/g3.119.400728 (2019).
- 18 5 Desta, Z. A. & Ortiz, R. Genomic selection: genome-wide prediction in plant improvement. *Trends in*
19 *plant science* 19, 592-601, doi:doi.org/10.1016/j.tplants.2014.05.006 (2014).
- 20 6 De Los Campos, G. *et al.* Predicting quantitative traits with regression models for dense molecular
21 markers and pedigree. *Genetics* 182, 375-385, doi:doi.org/10.1534/genetics.109.101501 (2009).
- 22 7 Li, Z. & Sillanpää, M. J. Overview of LASSO-related penalized regression methods for quantitative
23 trait mapping and genomic selection. *Theoretical and applied genetics* 125, 419-435,
24 doi:doi.org/10.1007/s00122-012-1892-9 (2012).
- 25 8 de los Campos, G., Vazquez, A. I., Fernando, R., Klimentidis, Y. C. & Sorensen, D. Prediction of
26 complex human traits using the genomic best linear unbiased predictor. *PLoS Genet* 9, e1003608,
27 doi:doi:10.1371/journal.pgen.1003608 (2013).
- 28 9 Meuwissen, T., Hayes, B. & Goddard, M. Prediction of total genetic value using genome-wide dense
29 marker maps. *Genetics* 157, 1819-1829 (2001).
- 30 10 Heffner, E. L., Jannink, J. L. & Sorrells, M. E. Genomic selection accuracy using multifamily
31 prediction models in a wheat breeding program. *The Plant Genome* 4, 65-75,
32 doi:doi.org/10.3835/plantgenome2010.12.0029 (2011).
- 33 11 Habier, D., Fernando, R. L., Kizilkaya, K. & Garrick, D. J. Extension of the Bayesian alphabet for
34 genomic selection. *BMC bioinformatics* 12, 186 (2011).
- 35 12 de los Campos, G., Hickey, J. M., Pong-Wong, R., Daetwyler, H. D. & Calus, M. P. Whole-genome
36 regression and prediction methods applied to plant and animal breeding. *Genetics* 193, 327-345,
37 doi:doi.org/10.1534/genetics.112.143313 (2013).
- 38 13 Dinesh, M. *et al.* Inheritance studies and validation of hybridity in guava (*Psidium guajava*). *Indian*
39 *Journal of Agricultural Sciences* 87, 42-45 (2017).
- 40 14 Pessanha, P. G. D. O. *et al.* Avaliação da Diversidade Genética em Acessos de *Psidium* spp. via
41 marcadores RAPD. *Revista Brasileira de Fruticultura* 33, 129-136, doi:doi.org/10.1590/s0100-
42 29452011000100018 (2011).
- 43 15 Silva, F. A. *et al.* Impact of Bayesian Inference on the Selection of *Psidium guajava*. *Scientific Reports*
44 10, 1-9, doi:doi.org/10.1038/s41598-020-58850-6 (2020).
- 45 16 Doyle, J. J. & Doyle, J. L. Isolation of plant DNA from fresh tissue. *Focus* 12, 39-40 (1990).
- 46 17 Guavamap, G. *Screening of microsatellite markers (SSRs) in Guava*
47 <<http://www.neiker.net/neiker/guavamap/for1-6a.htm>> (2008).
- 48 18 Pérez, P. & de Los Campos, G. Genome-wide regression and prediction with the BGLR statistical
49 package. *Genetics* 198, 483-495, doi:doi:10.1534/genetics.114.164442 (2014).

- 1 19 R Core Team (2018). R: A Language and Environment for Statistical Computing. R Foundation for
2 Statistical Computing, Vienna, Austria. URL:<http://www.R-project.org/>.
- 3 20 Spiegelhalter, D. J., Best, N. G., Carlin, B. P. & Van Der Linde, A. Bayesian measures of model
4 complexity and fit. *Journal of the royal statistical society: Series b (statistical methodology)* 64, 583-
5 639, doi:doi.org/10.1111/1467-9868.00353 (2002).
- 6 21 Wilberg, M. J. & Bence, J. R. Performance of deviance information criterion model selection in
7 statistical catch-at-age analysis. *Fisheries Research* 93, 212-221,
8 doi:doi.org/10.1016/j.fishres.2008.04.010 (2008).
- 9 22 Hoerl, A. E. & Kennard, R. W. Ridge regression: Biased estimation for nonorthogonal problems.
10 *Technometrics* 12, 55-67 (1970).
- 11 23 Tibshirani, R. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society:
12 Series B (Methodological)* 58, 267-288, doi:doi.org/10.1111/j.2517-6161.1996.tb02080.x (1996).
- 13 24 Heslot, N., Yang, H. P., Sorrells, M. E. & Jannink, J. L. Genomic selection in plant breeding: a
14 comparison of models. *Crop science* 52, 146-160, doi:doi:10.2135/cropsci2011.09.0297 (2012).
- 15 25 Viana, A. P. *et al.* Implementing genomic selection in sour passion fruit population. *Euphytica* 213,
16 228, doi:doi.org/10.1007/s10681-017-2020-3 (2017).
- 17 26 Crossa, J. *et al.* Genomic selection and prediction in plant breeding. *Journal of crop improvement*, 239,
18 doi:doi.org/10.1080/15427528.2011.558767 (2011).
- 19 27 Che, X. & Xu, S. Significance test and genome selection in bayesian shrinkage analysis. *International
20 Journal of Plant Genomics*, 1-11, doi:doi:10.1155/2010/893206 (2010).
- 21 28 Bihari, M. & Narayan, S. Genetic diversity, heritability, genetic advance and correlation coefficient in
22 guava (*Psidium guajava*). *Indian Journal of Agricultural Sciences* 81, 107-110 (2011).
- 23 29 Thaipong, K. & Boonprakob, U. Genetic and environmental variance components in guava fruit
24 qualities. *Scientia Horticulturae* 104, 37-47, doi:doi.org/10.1016/j.scienta.2004.07.008 (2005).

25

26 **Funding information**

27 This study was financed in part by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior - Brasil
28 (CAPES) - Finance Code 001. This study was financed by Fundação de Amparo à Pesquisa do Estado do Rio
29 de Janeiro (FAPERJ) – Finance Code E-26/010.001275/2015. This study was financed by Conselho Nacional
30 de Desenvolvimento Científico e Tecnológico (CNPq).

31 **Authors' contributions**

32 A.P.V supervision; F.A.S., E.A.S., J.A.V.S.O., J.D.G.A. and R.M.R. investigation; F.A.S. and C.C.G.C. writing
33 – original draft preparation; F.A.S. and L.S.G. formal analysis. All authors have reviewed the manuscript.

34 **Competing interests**

35 The authors declare no competing interests.

36 **Data Archiving Statement**

37 The full phenotypic information, breeding values, scripts and chains generated used in this study, have been
38 submitted at the *Open Science Framework* and was awarded the public doi identifier:

39 <https://doi.org/10.17605/OSF.IO/T8X7U>.

40

Figures

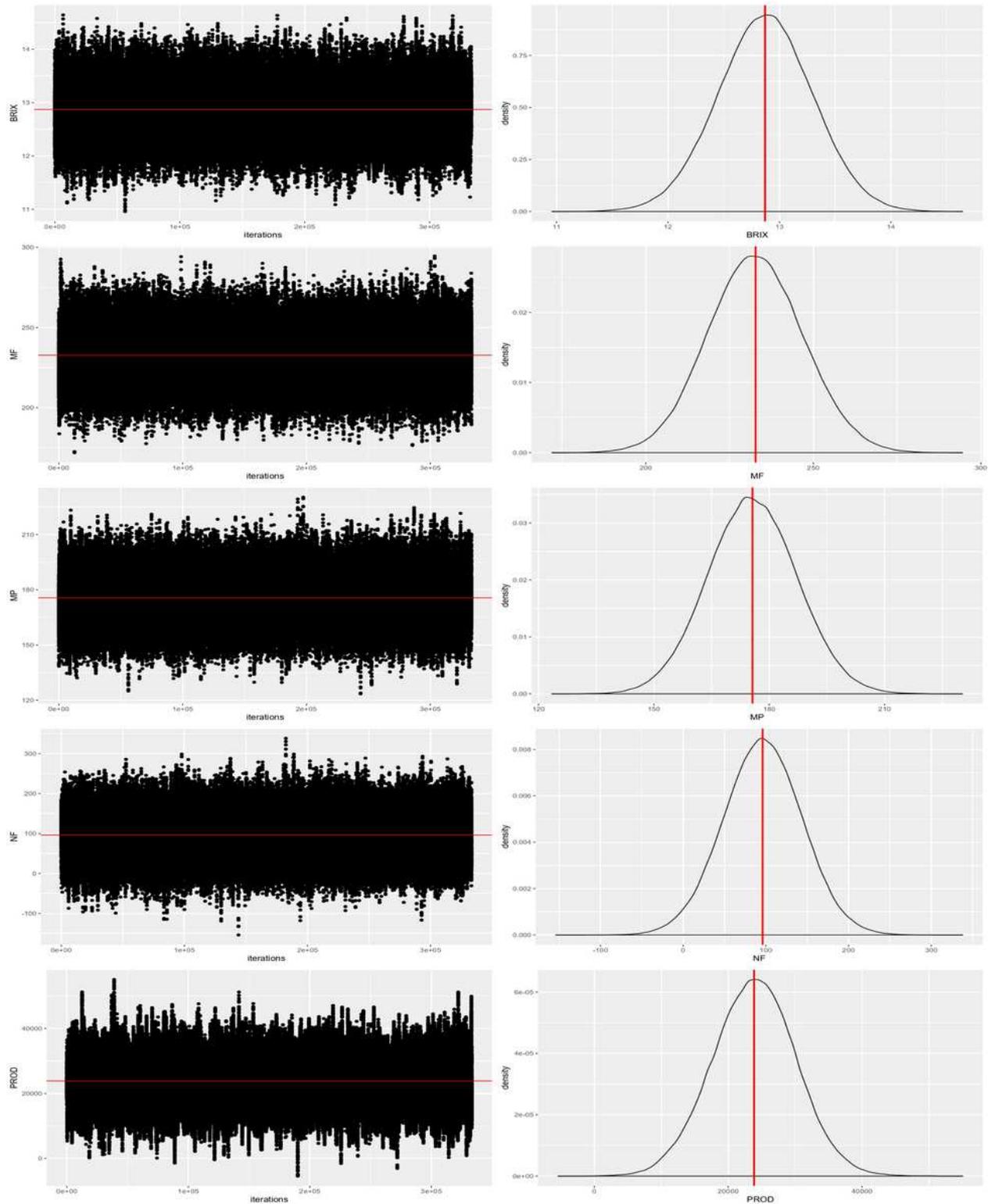


Figure 1

Markov and Monte Carlo chains with mean values (red line) and distribution curve for five variables observed in guava, generated to relate SSR marks to phenotypic observations.

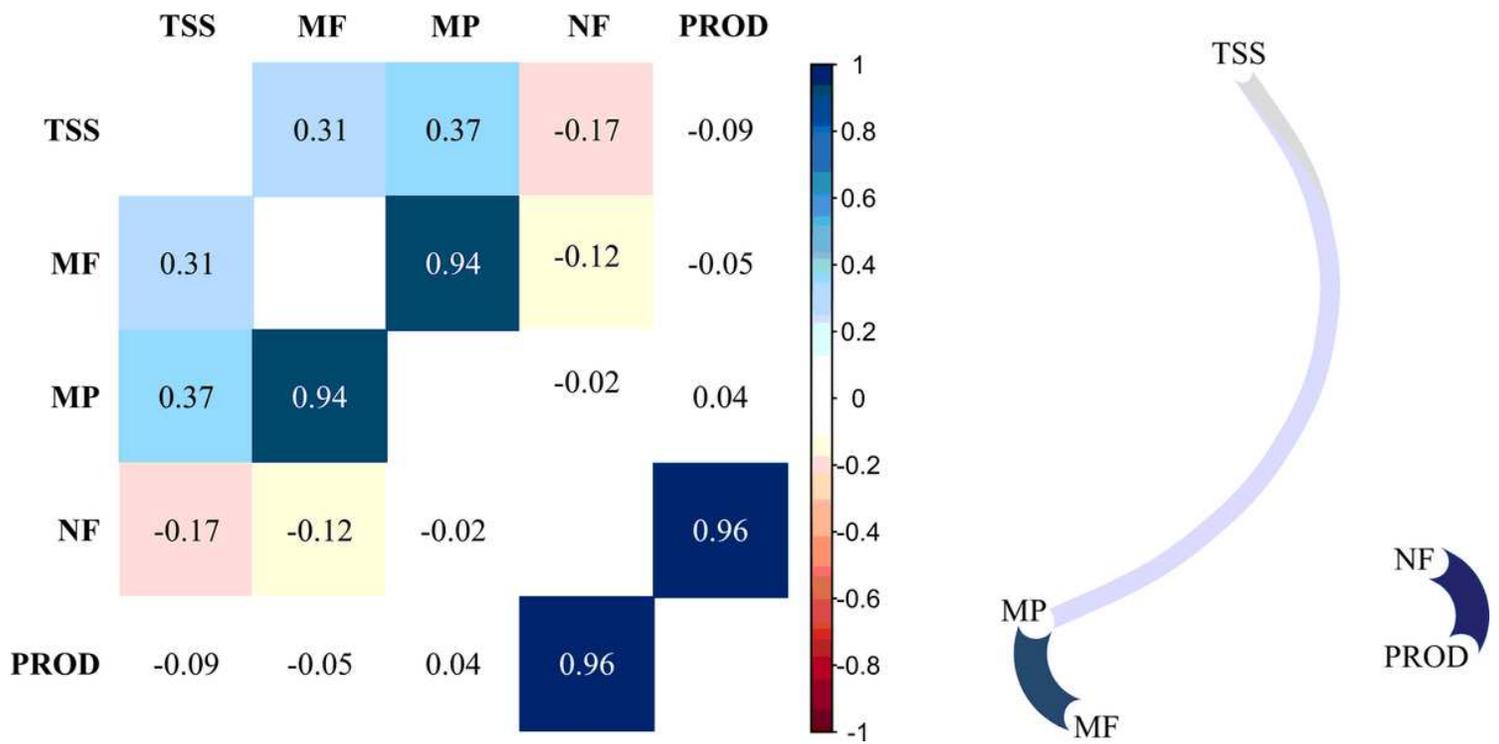


Figure 2

Genetic correlation between the soluble solids content (TSS), fruit mass (FM), pulp mass (PM), number of fruits per plant (NF), and production per plant (PROD) observed in guava (*Psidium guajava*), estimated using a model with SSR markers and Bayesian ridge regression - BRR.