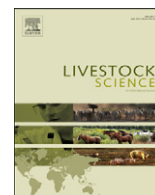


Contents lists available at [SciVerse ScienceDirect](http://SciVerse.Sciencedirect.com)

# Livestock Science

journal homepage: [www.elsevier.com/locate/livsci](http://www.elsevier.com/locate/livsci)

## Study of using marker assisted selection on a beef cattle breeding program by model comparison

F.M. Rezende<sup>a,\*</sup>, J.B.S. Ferraz<sup>a</sup>, J.P. Eler<sup>a</sup>, R.C.G. Silva<sup>a,b</sup>, E.C. Mattos<sup>a</sup>, N. Ibáñez-Escriche<sup>c</sup><sup>a</sup> Animal Breeding and Biotechnology Group, University of Sao Paulo, 13635-900 Pirassununga, SP, Brazil<sup>b</sup> Merial Igenity, 13091-908 Campinas, SP, Brazil<sup>c</sup> Genética i Millora Animal, IRTA-Cataluña, 25198 Lleida, Spain

### ARTICLE INFO

#### Article history:

Received 31 July 2011

Received in revised form

30 March 2012

Accepted 31 March 2012

#### Keywords:

Cross-validation

Genetic markers

SNP markers

Zebu cattle

### ABSTRACT

A data set of a commercial Nellore beef cattle selection program was used to compare breeding models that assumed or not markers effects to estimate the breeding values, when a reduced number of animals have phenotypic, genotypic and pedigree information available. This herd complete data set was composed of 83,404 animals measured for weaning weight (WW), post-weaning gain (PWG), scrotal circumference (SC) and muscle score (MS), corresponding to 116,652 animals in the relationship matrix. Single trait analyses were performed by MTDFREML software to estimate fixed and random effects solutions using this complete data. The additive effects estimated were assumed as the reference breeding values for those animals. The individual observed phenotype of each trait was adjusted for fixed and random effects solutions, except for direct additive effects. The adjusted phenotype composed of the additive and residual parts of observed phenotype was used as dependent variable for models' comparison. Among all measured animals of this herd, only 3160 animals were genotyped for 106 SNP markers. Three models were compared in terms of changes on animals' rank, global fit and predictive ability. Model 1 included only polygenic effects, model 2 included only markers effects and model 3 included both polygenic and markers effects. Bayesian inference via Markov chain Monte Carlo methods performed by TM software was used to analyze the data for model comparison. Two different priors were adopted for markers effects in models 2 and 3, the first prior assumed was a uniform distribution (U) and, as a second prior, was assumed that markers effects were distributed as normal (N). Higher rank correlation coefficients were observed for models 3\_U and 3\_N, indicating a greater similarity of these models animals' rank and the rank based on the reference breeding values. Model 3\_N presented a better global fit, as demonstrated by its low DIC. The best models in terms of predictive ability were models 1 and 3\_N. Differences due prior assumed to markers effects in models 2 and 3 could be attributed to the better ability of normal prior in handle with collinear effects. The models 2\_U and 2\_N presented the worst performance, indicating that this small set of markers should not be used to genetically evaluate animals with no data, since its predictive ability is restricted. In conclusion, model 3\_N presented a slight superiority when a reduce number of animals have phenotypic, genotypic and pedigree information. It could be attributed to the variation retained by markers and polygenic effects assumed together and the normal prior assumed to markers effects, that deals better with the collinearity between markers.

© 2012 Elsevier B.V. Open access under the [Elsevier OA license](http://www.elsevier.com/locate/elsevier-ol).

\* Corresponding author. Tel.: +55 19 35654105; fax: +55 19 35654107.

E-mail address: [frezende@usp.br](mailto:frezende@usp.br) (F.M. Rezende).

## 1. Introduction

Traditionally selection for economically relevant quantitative traits is realized based on two sources of information, pedigree and phenotype, under mixed model methodology that combines them to estimate the best linear unbiased prediction (BLUP; Henderson, 1984) of breeding value, which has been a successful approach.

Currently, a third source of information based on DNA markers became available. During last decades, genetic polymorphisms were described and their use on the selection was named marker assisted selection (MAS). The basic idea of MAS is to exploit statistical dependencies (linkage disequilibrium) existing in the joint distribution of markers and quantitative trait loci (QTL) genotypes (Gianola et al., 2003), which can be used to improve predictions of genetic merit of candidates for selection in a breeding program (Fernando and Grossman, 1989).

Developments of high dense single nucleotide polymorphism (SNP) genotyping have increased the interest for applying MAS at a genome wide scale, which has been termed genomic selection (Meuwissen et al., 2001). However, given the economic cost of this technology and the animal breeding schemes, its practical application remains mainly limited to dairy cattle, especially in the Holstein breed. Moreover, there is still not a consensus about what is the best methodology or strategy to apply it. The first methodology proposed to estimate markers effects was applying least squares method (Geldermann, 1975). In 1989, Fernando and Grossman estimated the effect of one locus under mixed models theory. Ridge regression methodology was also suggested to estimate markers effects, as a way to handle with dimensional problem and collinear effects (Whittaker et al., 2000). In the last decade, Bayesian penalized and non-parametric methods have been applied to markers effects estimations (De los Campos et al., 2009; Gianola et al., 2006; Habier et al., 2011; Meuwissen et al., 2001). In the meantime, several strategies to incorporate markers information on breeding programs have been proposed, as selection index, two and single step analyses and non-parametric methods (González-Recio et al., 2008; Legarra et al., 2009; Soller, 1978; VanRaden, 2008).

In beef cattle, the genomic selection have not been widely applied. As pointed out by Ibañez-Escriche and Gonzalez-Recio (2011) it could be due to several factors: the different sorts of organization in the breeding programs, the lack of systematic recording of phenotypic information, different breeding goals between populations and a smaller population size. In Zebu cattle, especially in Nellore breed, the most prevalent in the Brazilian beef industry, the number of studies on application of genetic markers is scarce. In fact, these studies have mainly been focused on investigation of associations between individual polymorphisms and quantitative traits (Aires et al., 2010; Ferraz et al., 2009; Pinto et al., 2010). To our knowledge, there are no scientific works reporting MAS to estimate the genomic breeding value with SNP markers on a Nellore cattle breeding program. The aim of the present study was carrying out a first evaluation of MAS implementation on a Nellore beef cattle breeding program,

applying different evaluation breeding models and comparing those approaches in terms of the changes on animals' rank, model global fit and predictive ability, when a reduced number of animals have phenotypic, genotypic and pedigree information available.

## 2. Materials and methods

### 2.1. Data

The data set for this study was obtained from the Agro-Pecuária CFM Ltda., a Brazilian beef cattle breeding company, whose genetic evaluation is carried out by the Animal Breeding and Biotechnology Group at the University of Sao Paulo ([www.usp.br/gmab](http://www.usp.br/gmab)). All evaluated animals were born between 1984 and 2009 and are progenies of bulls selected for production and reproduction traits under pasture conditions. The analyzed traits were weaning weight (WW), post-weaning gain (PWG), scrotal circumference (SC) and muscle score (MS). WW was recorded at around 205 days of age. Post-weaning gain was calculated as the weight gain between 205 and 550 days of age. Scrotal circumference measurements were taken at the greatest diameter of the scrotum, using a metal tape device, and were carried out at around 550 days of age. Muscle score was determined by a set of experts after visual evaluation of each animal at around 550 days of age. Animals with the poorest muscle content received a score of 1, and animals with the highest muscle content received a score of 6.

### 2.2. Data adjustment

Single trait analyses were performed by MTDFREML software (Boldman et al., 1995) to estimate fixed and random effects solutions for WW, PWG, SC and MS traits, under animal model, using the complete data set available for this herd. Fixed effects fitted were contemporary group composed by farm, year, season, sex and management group, age at measurement taken into account as a linear covariate and dam age as a quadratic covariate. Maternal additive and permanent environmental effects were fitted as random effects to WW and management group at weaning was fitted as a random effect for PWG, SC and MS.

Basic edits in the data set involved consistency checks of ages and measures, as well the elimination of records measured on animals with unknown dam, born from multiple sire group and on contemporary group composed of less than 5 animals. The descriptive statistics of the complete data set analyzed to estimate fixed and random effects solutions were presented in Table 1. The corresponding relationship matrix was composed of a total of 116,652 animals.

The individual records for all analyzed traits were adjusted for the same fixed and random effects fitted in the models used to estimate these effects solutions, except for animal effect. The adjusted phenotype represents the sum of direct additive and residual portions of the observed phenotype. The additive genetic effects were considered to be the reference breeding value for

**Table 1**

Summary statistics of the complete data set used to estimate fixed and random effects solutions for weaning weight (WW), post-weaning gain (PWG), scrotal circumference (SC) and muscle score (MS).

Trait	N	Mean	SD	CV
WW	83,404	188.93	27.35	14.48
PWG	68,424	114.47	32.25	28.18
SC	35,401	27.36	3.43	12.53
MS	63,854	3.61	1.00	27.84

N: number of data; SD: standard deviation; CV: coefficient of variation.

**Table 2**

Summary statistics for weaning weight (WW), post-weaning gain (PWG), scrotal circumference (SC) and muscle score (MS) of genotyped animals.

Trait	N	Mean	SD	CV
WW	3,042	207.05	22.90	11.06
PWG	3,033	125.19	30.50	24.36
SC	2,664	27.18	3.36	12.36
MS	3,149	3.58	1.17	32.68

N: number of data; SD: standard deviation; CV: coefficient of variation.

genotyped and non-genotyped animals. Adjusted phenotype and additive genetic values will be used in the next steps for model comparison.

### 2.3. Genotypic data

From this herd, a sample of 3549 animals were genotyped, of which 3160 composed the genotypic data set, once they had known sire and dam. The genotypic data set consisted of 377 females and 2783 males. A total of 3010 dams and 752 sires were represented on this data set, of these 104 dams and 278 sires were also genotyped. Most of dams have just one progeny analyzed and 46 sires have more than 10 progenies on genotypic data set, from that 21 were genotyped sires. Descriptive statistics for WW, PWG, SC and MS of genotyped animals are presented in Table 2.

Animals were genotyped to 222 SNP markers identified on *Bos taurus* breeds, represented by 123 SNP markers described on literature to be associated with genes that affect production traits expression, of which 85–90% are on transcript regions and 10–15% are on promote regions, and 99 SNP markers used in paternity tests, which, although were not described to be associated with any biological function, are distributed along the genome and, then, they could be in linkage disequilibrium with genes that are affecting these traits. The genotyping process was carried out in laboratories, licensed by Merial/Igenity®, a company that holds the rights of use of those markers.

Allelic and genotypic frequencies for each marker were estimate by simple count of different alleles and genotypes, using PROC FREQ from SAS. SNP markers with minor allelic frequency lower than 5% were removed from the data set, and, after that, 106 genetic markers were kept for analyses. This reduction of the number of markers kept for analyses compared to genotyped markers is

due to the fact that those markers were described on *Bos taurus* breeds and Nellore is a *Bos indicus* breed resulting in many markers fixed or presenting a minor allelic frequency less than 5%.

### 2.4. Statistical models

Three linear mixed models were used to predict genetic and markers effects for genotyped animals, that represent less than 8% of all animals measured in this herd.

*Model 1:* This model included only polygenic effects and can be expressed, in matrix algebra notation, as

$$\mathbf{y} = \mu + \mathbf{Za} + \mathbf{e}$$

where  $\mathbf{y}$  is a vector  $n \times 1$  of the adjusted phenotype ( $n$ =number of records);  $\mu$  is the overall mean;  $\mathbf{a}$  is a vector  $q \times 1$  of additive genetic polygenic effects ( $q$ =number of animals on relationship matrix);  $\mathbf{Z}$  is the additive genetic effects design matrix of order  $n \times q$ ; and  $\mathbf{e}$  is the residuals vector.

*Model 2:* This model included only markers effects and can be described as

$$\mathbf{y} = \mu + \mathbf{Xg} + \mathbf{e}$$

where  $\mathbf{y}$  is a vector  $n \times 1$  of adjusted phenotype ( $n$ =number of records);  $\mu$  is the overall mean;  $\mathbf{g}$  is a vector  $p \times 1$  of allele substitution marker effect ( $p$ =number of analyzed markers);  $\mathbf{X}$  is the incidence genotype matrix of order  $n \times p$ , whose elements were set up as an additive model, with values 1, 2 or 3 for aa, Aa and AA, respectively; and  $\mathbf{e}$  is the residuals vector.

*Model 3:* This model included both markers and polygenic effects

$$\mathbf{y} = \mu + \mathbf{Xg} + \mathbf{Za} + \mathbf{e}$$

where  $\mathbf{y}$  is a vector  $n \times 1$  of adjusted phenotype ( $n$ =number of records);  $\mu$  is the overall mean;  $\mathbf{g}$  is a vector of  $p \times 1$  allele substitution marker effect ( $p$ =number of analyzed markers);  $\mathbf{X}$  is the incidence genotype matrix of order  $n \times p$ , whose elements were set up as an additive model, with values 1, 2 or 3 for aa, Aa and AA, respectively;  $\mathbf{a}$  is a vector  $q \times 1$  of additive genetic polygenic effects ( $q$ =number of animals on relationship matrix);  $\mathbf{Z}$  is the additive genetic effects design matrix of order  $n \times q$ ; and  $\mathbf{e}$  is the residuals vector.

The number of records corresponds to the number of genotyped animals for each analyzed trait, as it was presented in Table 2. The relationship matrix considered in models 1, 3\_U and 3\_N was the same previous described, composed of 116,652 animals.

### 2.5. Statistical analyses

Bayesian inference via Markov chain Monte Carlo methods was used to analyze the data. The distribution of data was assumed to be a normal distribution

$$\mathbf{y} | \theta_j \sim N(\mu_j, \mathbf{I}\sigma_e^2)$$

where the subscript  $j=1, 2$  and  $3$  denotes model 1, model 2 and model 3, respectively;  $\theta_j$  are the unknown parameters

for each model;  $\sigma_e^2$  is the residual variance. Note that,  $\mu_1 = \mu + \mathbf{Za}$ ,  $\mu_2 = \mu + \mathbf{Xg}$  and  $\mu_3 = \mu + \mathbf{Xg} + \mathbf{Za}$ .

For all models, a uniform prior was used for  $\mu$ , and residual effects were assumed to be distributed as  $N(0, \mathbf{I}\sigma_e^2)$  with  $\sigma_e^2$  distributed as  $\text{Inv-}\chi^2(v, s^2)$ , with  $v = -2$  and  $s^2 = 0$ . In models 1 and 3, additive genetic effects were assumed to be distributed as  $N(0, \mathbf{A}\sigma_a^2)$ , where  $\mathbf{A}$  is the relationship matrix and  $\sigma_a^2$  is the additive variance assumed to be distributed as  $\text{Inv-}\chi^2(v, s^2)$ , with  $v = -2$  and  $s^2 = 0$ . Two different priors were used for markers effects in models 2 and 3. The first prior assumed was a uniform distribution, which is similar to the regression analysis described by Meuwissen et al. (2001) where markers effects were estimate by least squares method. The second prior assumed that markers effects were distributed as  $N(0, \mathbf{I}\sigma_g^2)$ , where  $\sigma_g^2$  is the marker additive variance assumed to be distributed as  $\text{Inv-}\chi^2(v, s^2)$ , with  $v = -2$  and  $s^2 = 0$ . This approach, called as Bayesian ridge regression, was proposed by Gianola et al. (2003), it is similar to the ridge regression estimator proposed by Whittaker et al. (2000) in a frequentist context. Next, models 2 and 3 will be identified according to the prior assumed to markers effects as 2\_U and 3\_U for uniform distribution prior and as 2\_N and 3\_N for a normal distribution prior. Therefore, model comparison will be considered to have five different models. The fully conditional posterior distributions for all models were normal distribution for  $\mu$ ,  $\mathbf{g}$  and  $\mathbf{a}$  and inverted scaled chi-square distribution for  $\sigma_a^2$ ,  $\sigma_g^2$  and  $\sigma_e^2$ .

The analyses to estimate markers and genetic effects on models 1, 2\_U, 2\_N, 3\_U and 3\_N were performed by TM software (Legarra et al., 2008) modified to include markers effects. For each analysis, a single chain with a total of 1,000,000 Gibbs sampler iterations, a burn-in period of 80,000 and a thin interval of 100 samples were used. Convergence was tested separately for all unknown parameters using the Raftery and Lewis (1992) algorithm, the Z criterion of Geweke (1992), the Monte Carlo sampling errors computed using the time-series procedures described by Geyer (1992) and a visual check of the chain plot. All those tests were performed by BOA (Bayesian Output Analysis) diagnostic program from R package.

## 2.6. Model comparison

The animal's rank based on the reference breeding values estimated using the complete data set of this herd was compared to the rank given by the breeding values estimated on models 1, 2\_U, 2\_N, 3\_U and 3\_N. The association between the rank of genotyped animals based on reference breeding values with the rank based on breeding values estimated by the five models were calculated using two non-parametric measures of ordinal association computed by PROC CORR from SAS.

*Spearman's rank correlation coefficient* ( $\rho_s$ ). Spearman's rank correlation coefficient is computed by ranking the data and using the ranks in the Pearson product-moment correlation formula, as shown below. In case of ties, averaged ranks are used.

$$\rho_s = \frac{\sum_i((R_i - \bar{R})(S_i - \bar{S}))}{\sqrt{\sum_i(R_i - \bar{R})^2 \sum_i(S_i - \bar{S})^2}}$$

where  $R_i$  is the rank of each observation based on reference breeding value;  $S_i$  is the rank of each observation based on genetic value estimated from models 1, 2\_U, 2\_N, 3\_U or 3\_N;  $\bar{R}$  is the mean of  $R_i$  values;  $\bar{S}$  is the mean of  $S_i$  values.

*Kendall's tau-b rank correlation coefficient* ( $\tau_b$ ). The data are double sorted, first the observations are ranked according to the reference breeding values and, then, re-ranked according to the breeding values obtained from models 1, 2\_U, 2\_N, 3\_U or 3\_N. Kendall's tau-b is calculated based on the number of concordant and discordant pairs of observations, as

$$\tau_b = \frac{n_c - n_d}{(1/2)n(n-1)}$$

where  $n_c$  is the number of concordant pairs;  $n_d$  is the number of discordant pairs;  $n$  is the total number of measured animal.

The models 1, 2\_U, 2\_N, 3\_U and 3\_N were also compared in terms of global fit and predictive ability by the following methods.

*Deviance information criteria* (DIC). The DIC compares the global adjusted quality of two or more models, accounting for model complexity (Spiegelhalter et al., 2002). For a particular model M, the DIC is defined as

$$\text{DIC} = 2\bar{D} - D(\bar{\theta}_M)$$

The term  $\bar{D} = -2 \int [\log p(y|\theta_M)]p(\theta_M|y, M)d\theta_M = E_{\theta_M|y}[D(\theta_M)]$  is the posterior expectation of the deviance  $D(\theta_M)$ , and  $D(\bar{\theta}_M) = -2 \log p(y|\bar{\theta}_M)$  is the deviance evaluated at the posterior mean of the parameter vector  $\theta_M$ . DIC expression is the result of combining both terms, where  $\bar{D}$  is a measure of model fit and  $\bar{D} - D(\bar{\theta}_M)$  is related to the effective number of parameters ( $p_d$ ). Models with smaller DIC exhibit a better global fit after accounting for model complexity. Differences in DIC of more than 7 are considered as important by Spiegelhalter et al. (2002).

*K-fold cross-validation*. K-fold cross-validation approach was used to evaluate the models based on their ability to predict data. The entire genotypic data set was split into a training set for model fitting and a validation set to test the predictive ability of the model, using two distinct strategies: (1) 1-fold earlier progeny cross-validation, training set was composed of the older animals and validation set was composed of the youngest individuals, born in 2009, representing around of 25% of genotyped animals. This temporal strategy of splitting the data set has a strong relation with breeding programs, where the phenotype of ancestors and the relationship information are used to estimate the breeding values of the youngest animals; (2) 4-fold cross-validation, partitioning genotypic data set into four disjoint subsets, each with approximately one-fourth of the records, by taking random samples of data points. The cross-validation procedure used three of the four subsets for training set and the remaining subset was used for validation set. This procedure was realized four times for changing the subset used to test the predictive ability of the model. This non-temporal strategy permitted to evaluate the predictive ability of models when one-fourth of the records, defined randomly, were not measured.

Two different criteria were used to compare the predictive ability of the models.

**Mean squared error.** The mean squared error (MSE) was computed as

$$\text{MSE} = \frac{1}{n_{\text{data}}} \sum_{i=1}^{n_{\text{data}}} (y - \hat{y})^2$$

where  $y$  and  $\hat{y}$  correspond to the observed and predicted observations, respectively; and  $n_{\text{data}}$  is the number of data points in the validation subset. Models having the smallest MSE were regarded as those with the best predictive ability.

**Pearson's correlation.** Pearson's correlation ( $\rho$ ) between observed and predicted observations was calculated as

$$\rho_{y,\hat{y}} = \frac{1}{n_{\text{data}}} \sum_{i=1}^{n_{\text{data}}} \frac{\text{cov}(y,\hat{y})}{\sigma_y \sigma_{\hat{y}}}$$

where  $\text{cov}(y,\hat{y})$  is the covariance estimate between observed and predicted records;  $\sigma_y$  and  $\sigma_{\hat{y}}$  are the estimates of standard deviations of observed and predicted records; and, as above,  $n_{\text{data}}$  is the number of data points in the validation subset. The model providing the highest correlation was considered as the one with the best predictive ability.

### 3. Results

Posterior distributions convergence for all unknown parameters in the five models tested by the Raftery and Lewis method, Z criterion of Geweke and a visual check of the chain plot did not detect any lack of convergence. The Monte Carlo sampling error for posterior distributions of all unknown parameters was irrelevant, as it was at least 20 times lower than the standard deviation of posterior distributions. This low magnitude of Monte Carlo sampling error also suggests convergence of chains, according to [Blasco et al. \(2003\)](#).

Spearman's ( $\rho_s$ ) and Kendall's tau-b ( $\tau_b$ ) rank correlation coefficients for WW, PWG, SC and MS, calculated to compare the ranks of animals based on breeding values estimated by models 1, 2\_U, 2\_N, 3\_U and 3\_N with the rank based on reference breeding values estimated using this herd complete data set, are given in [Table 3](#). In general,  $\tau_b$  coefficients were lower than  $\rho_s$  coefficients for all models and traits, suggesting  $\tau_b$  as a more conservative

measure. The lowest values for both measures of association were exhibited by models 2\_U and 2\_N, indicating greatest divergences in the animals' rank based on these models' estimates of breeding values and the rank based on the reference breeding values. The highest values for models 3\_U and 3\_N suggest small differences between the rank of animals by these models and the rank given by the reference breeding values. Furthermore, models that assumed normal prior distribution to markers effects presented rank correlation coefficients slightly higher.

The posterior means and the highest 95% posterior density intervals (HPD95%) of additive, markers and residual variances and coefficients of heritability estimated for WW, PWG, SC and MS by models 1, 2\_U, 2\_N, 3\_U and 3\_N are presented in [Table 4](#). Posterior distributions of the additive variances and heritability for models 1, 3\_U and 3\_N within traits were similar and their HPD95% were overlapped. Comparing the estimates of residual variances, the lowest values within traits were observed for model 3\_N and the highest values for models 2\_U and 2\_N.

[Table 5](#) shows the estimates of the deviance information criterion (DIC), the deviance ( $\bar{D}$ ) and the effective number of parameters ( $p_d$ ) for all models. The DIC criterion is based on a balance between the fit of the data to the model and the corresponding complexity of the model and its behavior is based on the ability to make short-term predictions of a repeat set of similar data. In our case, model 3\_N presented the highest global fit for all analyzed traits, except for WW.

Two cross-validation strategies, 1-fold and 4-fold, and two criteria, mean squared error (MSE) and Pearson's correlation, were adopted to evaluate the predictive ability of models 1, 2\_U, 2\_N, 3\_U and 3\_N. [Tables 6 and 7](#) present the results of MSE and Pearson's correlation for 1-fold earlier progeny cross-validation analyses, respectively. In terms of MSE, smaller values were exhibited by model 1 for WW and PWG and by model 3\_N for SC and MS, suggesting a better predictive ability of these models for those traits. Estimates of Pearson's correlation of models 1 and 3\_N were the same for WW and PWG and, for SC and MS, model 3\_N presented values slightly higher than model 1.

In [Tables 8 and 9](#) are presented the results of MSE and Pearson's correlation for 4-fold cross-validation analyses, respectively. In agreement with 1-fold results, the best

**Table 3**

Estimates of Spearman's ( $\rho_s$ ) and Kendall's tau-b ( $\tau_b$ ) rank correlation coefficients from models 1, 2\_U, 2\_N, 3\_U and 3\_N compared with the rank based on the reference breeding values.

Correlation	Trait	Model 1	Model 2_U	Model 2_N	Model 3_U	Model 3_N
$\rho_s$	WW	0.47	0.22	0.25	0.66	0.82
	PWG	0.53	0.36	0.37	0.83	0.90
	SC	0.57	0.31	0.30	0.94	0.97
	MS	0.57	0.29	0.29	0.84	0.89
$\tau_b$	WW	0.32	0.15	0.17	0.48	0.63
	PWG	0.37	0.24	0.25	0.64	0.72
	SC	0.40	0.21	0.20	0.80	0.85
	MS	0.41	0.20	0.20	0.65	0.71

WW: weaning weight; PWG: post-weaning gain; SC: scrotal circumference; MS: muscle score; U: uniform distribution prior for markers effects; N: normal distribution prior for markers effects.

**Table 4**

Estimates of posterior means and 95% highest posterior density intervals in parenthesis for additive ( $\sigma_a^2$ ), markers ( $\sigma_g^2$ ) and residual ( $\sigma_e^2$ ) variances and heritabilities ( $h^2$ ) of weaning weight (WW), post-weaning gain (PWG), scrotal circumference (SC) and muscle score (MS).

Model	Parameter	PDES	GPDES	PE	EM
1	$\sigma_a^2$	25.47 (15.10, 37.28)	127.01 (92.80, 160.11)	3.66 (2.85, 4.43)	0.38 (0.28, 0.50)
	$\sigma_e^2$	133.21 (121.33, 144.11)	259.34 (229.36, 289.12)	3.31 (2.69, 3.93)	0.61 (0.51, 0.69)
	$h^2$	0.16 (0.09, 0.23)	0.33 (0.25, 0.41)	0.52 (0.43, 0.62)	0.39 (0.29, 0.48)
2_U	$\sigma_e^2$	156.24 (148.53, 164.46)	367.53 (349.20, 386.36)	6.47 (6.12, 6.83)	0.93 (0.89, 0.98)
2_N	$\sigma_g^2$	0.02 (0.00, 0.05)	0.50 (0.25, 0.76)	0.01 (0.005, 0.02)	0.001 (0.0004, 0.002)
	$\sigma_e^2$	156.24 (148.22, 164.06)	367.44 (348.73, 386.07)	6.50 (6.16, 6.88)	0.93 (0.89, 0.98)
3_U	$\sigma_a^2$	25.35 (14.67, 36.82)	112.41 (79.13, 148.18)	3.48 (2.72, 4.33)	0.36 (0.25, 0.46)
	$\sigma_e^2$	133.97 (122.16, 145.04)	264.68 (233.21, 294.66)	3.31 (2.67, 3.94)	0.62 (0.52, 0.71)
	$h^2$	0.16 (0.09, 0.23)	0.30 (0.21, 0.38)	0.51 (0.41, 0.61)	0.36 (0.26, 0.47)
3_N	$\sigma_a^2$	25.50 (15.22, 35.46)	120.98 (87.21, 153.11)	3.55 (2.78, 4.40)	0.37 (0.27, 0.48)
	$\sigma_g^2$	0.007 (0.00, 0.03)	0.24 (0.07, 0.44)	0.005 (0.00, 0.01)	0.0006 (0.00, 0.001)
	$\sigma_e^2$	133.09 (122.03, 143.62)	257.54 (228.98, 287.30)	3.28 (2.66, 3.91)	0.60 (0.51, 0.69)
	$h^2$	0.16 (0.10, 0.22)	0.32 (0.24, 0.40)	0.52 (0.42, 0.62)	0.38 (0.28, 0.48)

U: uniform distribution prior for markers effects; N: normal distribution prior for markers effects.

**Table 5**

Estimates of the deviance information criterion (DIC), the deviance ( $\bar{D}$ ) and the effective number of parameters ( $p_d$ ) from models 1, 2\_U, 2\_N, 3\_U and 3\_N.

Model	Parameter	WW	PWG	SC	MS
1	DIC	23,934.89	26,316.04	11,965.48	8,408.07
	$\bar{D}$	23,507.92	25,457.25	10,738.52	7,359.02
	$p_d$	426.97	858.80	1226.96	1,049.04
2_U	DIC	24,105.70	26,626.99	12,640.97	8,820.40
	$\bar{D}$	23,997.42	26,518.86	12,532.96	8,712.36
	$p_d$	108.28	108.12	108.01	108.05
2_N	DIC	24,005.91	26,573.64	12,598.50	8,766.28
	$\bar{D}$	23,995.46	26,518.48	12,544.31	8,716.90
	$p_d$	10.45	55.16	54.20	49.38
3_U	DIC	24,037.92	26,377.91	11,984.66	8,455.03
	$\bar{D}$	23,526.44	25,521.65	10,731.66	7,395.34
	$p_d$	511.47	856.26	1253.00	1059.69
3_N	DIC	23,936.58	26,297.91	11,941.03	8,390.72
	$\bar{D}$	23,505.37	25,437.93	10,709.36	7,337.42
	$p_d$	431.21	859.98	1231.67	1,053.30

WW: weaning weight; PWG: post-weaning gain; SC: scrotal circumference; MS: muscle score; U: uniform distribution prior for markers effects; N: normal distribution prior for markers effects.

models in terms of predictive ability were models 1 and 3\_N, with minor differences between them. For WW, model 1 presented the lower mean value of MSE and the highest mean value for Pearson's correlation, therefore, a better predictive ability of this model for that trait. Model 3\_N was the preferred model for PWG, SC and MS, as indicated by its lowest values for MSE and highest values for Pearson's correlation.

Differences of predictive ability between model 1 and model 3\_N were very small for both cross-validation strategies and criteria. Nevertheless, the worst model in

**Table 6**

Estimates of mean squared error from models 1, 2\_U, 2\_N, 3\_U and 3\_N by 1-fold earlier progeny cross-validation.

Trait	Model 1	Model 2_U	Model 2_N	Model 3_U	Model 3_N
WW	148.50	158.13	149.33	157.76	148.71
PWG	362.05	390.62	369.14	390.13	363.85
SC	6.06	6.72	6.35	6.36	6.03
MS	1.18	1.23	1.20	1.20	1.17

WW: weaning weight; PWG: post-weaning gain; SC: scrotal circumference; MS: muscle score; U: uniform distribution prior for markers effects; N: normal distribution prior for markers effects.

**Table 7**

Estimates of Pearson's correlation from models 1, 2, 3 by 1-fold earlier progeny cross-validation.

Trait	Model 1	Model 2_U	Model 2_N	Model 3_U	Model 3_N
WW	0.06	0.00	-0.01	0.02	0.06
PWG	0.11	0.04	0.06	0.07	0.11
SC	0.22	0.09	0.10	0.20	0.24
MS	0.14	0.08	0.09	0.13	0.15

WW: weaning weight; PWG: post-weaning gain; SC: scrotal circumference; MS: muscle score; U: uniform distribution prior for markers effects; N: normal distribution prior for markers effects.

terms of predictive ability was model 2\_U, that presented the highest values of MSE and the lower values of Pearson's correlation for both cross-validation strategies.

#### 4. Discussion

The advantage of including the information of a small set of molecular markers, when a reduced number of animals have phenotypic and genealogy information available, has been evaluated here by the comparison of five different models. First, we evaluated changes in the rank of animals based on reference breeding values estimated using the complete data set and the rank based

**Table 8**

Estimates of mean squared error from models 1, 2\_U, 2\_N, 3\_U and 3\_N to subsets 1, 2, 3 and 4, by 4-fold cross-validation.

Trait	Subset	Model 1	Model 2_U	Model 2_N	Model 3_U	Model 3_N
WW	1	152.58	162.39	155.22	159.27	152.57
	2	147.51	155.77	148.61	155.39	147.56
	3	155.01	169.12	158.28	167.72	155.44
	4	160.25	175.33	164.35	172.58	160.68
	Mean	153.84	165.65	156.62	163.74	154.06
PWG	1	389.62	407.87	403.55	399.50	386.51
	2	323.37	344.98	327.42	339.77	322.81
	3	352.40	386.80	371.39	369.19	351.29
	4	380.73	401.76	396.74	385.47	376.95
	Mean	361.53	385.35	374.78	373.48	359.39
SC	1	7.21	7.62	7.57	7.27	7.15
	2	5.41	6.00	5.85	5.57	5.38
	3	5.98	6.66	6.42	6.23	5.96
	4	6.18	7.04	6.75	6.42	6.16
	Mean	6.20	6.83	6.65	6.37	6.16
MS	1	1.01	1.03	1.01	1.02	1.00
	2	0.87	0.92	0.89	0.89	0.86
	3	0.85	0.89	0.87	0.85	0.84
	4	1.00	1.07	1.03	1.06	1.00
	Mean	0.93	0.98	0.95	0.96	0.93

WW: weaning weight; PWG: post-weaning gain; SC: scrotal circumference; MS: muscle score; U: uniform distribution prior for markers effects; N: normal distribution prior for markers effects.

**Table 9**

Estimates of Pearson's correlation from models 1, 2\_U, 2\_N, 3\_U and 3\_N to subsets 1, 2, 3 and 4, by 4-fold cross-validation.

Trait	Subset	Model 1	Model 2_U	Model 2_N	Model 3_U	Model 3_N
WW	1	0.12	0.00	-0.02	0.06	0.12
	2	0.10	0.01	0.01	0.04	0.10
	3	0.14	-0.02	-0.03	0.02	0.12
	4	0.16	-0.03	-0.02	0.02	0.15
	Mean	0.13	-0.01	-0.02	0.04	0.12
PWG	1	0.25	0.18	0.19	0.24	0.27
	2	0.17	0.09	0.09	0.15	0.17
	3	0.26	0.10	0.13	0.20	0.26
	4	0.26	0.17	0.17	0.24	0.28
	Mean	0.24	0.14	0.15	0.21	0.25
SC	1	0.27	0.18	0.18	0.27	0.29
	2	0.30	0.16	0.14	0.29	0.31
	3	0.28	0.12	0.12	0.25	0.29
	4	0.31	0.12	0.12	0.27	0.31
	Mean	0.29	0.15	0.14	0.27	0.30
MS	1	0.16	0.12	0.12	0.17	0.17
	2	0.19	0.11	0.12	0.18	0.21
	3	0.22	0.13	0.13	0.22	0.24
	4	0.18	0.07	0.06	0.13	0.17
	Mean	0.19	0.11	0.11	0.18	0.20

WW: weaning weight; PWG: post-weaning gain; SC: scrotal circumference; MS: muscle score; U: uniform distribution prior for markers effects; N: normal distribution prior for markers effects.

on breeding values estimated by each model on a reduced data set. Later, models were compared regarding the overall fit and predictive ability.

Higher divergences on animals' rank were detected by  $\tau_b$  than by  $\rho_s$  for all models and traits, they are in

agreement with Kendall (1947) that described for a large number of data,  $\tau_b$  is about two-thirds of the value of  $\rho_s$ . However, the interpretations of both association measures, Spearman's ( $\rho_s$ ) and Kendall's tau-b ( $\tau_b$ ) rank correlation coefficients, were very similar and led to the

same inferences. The same pattern of rank correlation coefficients was observed for all analyzed traits. Higher divergences on the rank of animals were exhibited by models 2\_U and 2\_N, as indicated by their low correlations with the rank of animals given by reference breeding values. This divergence can be due to the small set of markers analyzed, that were not enough to explain all additive effect contained in the adjusted phenotype. Nevertheless, model 1 presented intermediate values of rank correlation coefficients, suggesting that this model was more capable to retain the additive part of adjusted phenotype than models 2\_U and 2\_N. It was expected considering the information provided by the dense relationship matrix available for those animals. Furthermore, the inclusion of markers and polygenic effects together in models 3\_U and 3\_N lead to greater resemblance between the rank of animals based on these models' estimates of breeding values and the rank using the reference breeding values. The combination of both effects enhanced these models' efficiency to retain the additive effect contained on adjusted phenotype, when a reduced data set is available. The prior assumed for markers effects slightly affected the estimations of direct additive genetic effects between least square and ridge regression methods. The higher rank correlation coefficients observed in models that assumed normal prior distribution to markers effects can be attributed to their better ability to handle collinear effects (Whittaker et al., 2000).

Posterior means of coefficients of heritability estimated by models 1, 3\_U and 3\_N within traits were in agreement with the estimates obtained using this herd complete data set, that were 0.19, 0.21, 0.50 and 0.21 for WW, PWG, SC and MS, respectively. Since the HPD95% of those models were overlapped, their estimates of heritability cannot be considered different. Posterior means of heritability for WW, PWG, MS and SC obtained in the present study for all models were in the same range of estimates published in the last genetic evaluation of this herd, which were 0.22, 0.32, 0.45 and 0.25, respectively (personal communication). Horimoto et al. (2007) and Van Melis et al. (2010) while analyzing other samples of the same Nellore population reported estimations of heritability for WW that were out of the uncertainty interval presented here for this trait, 0.28 and 0.55, respectively. The authors described estimates in the same range of the ones observed here for PWG (0.32 and 0.25) and SC (0.55 and 0.42). And, for MS their estimations were also out of uncertainty interval presented here, 0.16 and 0.23.

The comparison of the residual variance estimated by each model showed that model 3\_N presented the lowest values for all analyzed traits, although small differences were observed between the estimates of residual variance of models 1, 3\_U and 3\_N. These results indicated their greater capacity for explaining the data variation in comparison with models 2\_U and 2\_N. Almost no difference was observed due to the prior assumed to markers effects on models 2\_U and 2\_N, that can be attributed to the absence of dimensional problem on the estimation of markers effects, since the number of markers was not greater than the number of observations. The problem of dimension was discussed by Meuwissen et al. (2001) as

the major limitation to the use of least squares method, since it leads to biased estimations of markers effects. Although there is no problem of dimension on the analyzed data set, the assumption of a uniform prior on model 3\_U led to a slim confounding on the estimations of polygenic and markers effects that is observed comparing the residual and additive variances of models 1 and 3\_U. It was expected that the inclusion of markers effects on model 3\_U should reduce the residual variance regarding model 1 without changes on the estimates of additive variance, but it was not observed here. The slight superiority of model 3\_N could be attributed to the variation retained by markers and polygenic effects assumed together and the normal prior assumed to markers effects, that deals better with the collinearity between markers, as was discussed by Whittaker et al. (2000).

The different criteria adopted to compare the five models provided analogous results. In terms of global fit, the deviance information criterion (DIC) clearly favored model 3\_N, except for WW, that model 1 presented the lowest value. Once DIC differences below 7 were not important (Spiegelhalter et al., 2002), those models presented the same global fit quality for WW.

The predictive ability of the models evaluated by mean squared error (MSE) and Pearson's correlation was consistent across *k*-fold cross-validation approach. Models 1 and 3\_N were equivalent in terms of predictive ability for all analyzed traits, with slight superiority of model 1 for WW and of model 3\_N for SC, in both cross-validations results. Differences between models 3\_U and 3\_N due to the prior assumed to markers effects became evident, with a worse predictive ability of model 3\_U. This can be attributed to the difficulty of model 3\_U in dealing with the existence of collinearity between markers effects, which do not occur with the model 3\_N. Given the small SNP panel available to this Nellore population, it was demonstrated that the predictive ability of the next generation performance or of randomly removed data based only on markers information (models 2\_U and 2\_N) was very limited. These results are in agreement with Lande and Thompson (1990) that reported that the selection on marker loci alone is more efficient when the proportion of the additive genetic variance explained by the marker loci exceeds the heritability of the character, which was not observed in this study.

Even though almost no differences on the predictive ability were observed when markers and polygenic effects were assumed together, the great advantage of including both effects was demonstrated by the rank correlation coefficients. The ranks given by breeding values estimates of models 3\_U and 3\_N were very similar with the rank based on the reference breeding values. Finally, the use of a normal prior for markers effects instead of a uniform prior is justified by the ability of this approach in handling with collinear effects between SNP markers.

## 5. Conclusion

In terms of rank correlation, a slight advantage was observed in the models that assumed polygenic and markers effects together. The highest global fit was presented



by the model that assumed both effects and a normal prior to markers effects (model 3\_N). No differences of predictive ability were observed between model 1, which included only polygenic effects, and model 3\_N. The worst overall fit and predictive ability were exhibited when only markers effects were considered in the model, regardless of prior assumed to these effects. In conclusion, for a reduced data set the inclusion of polygenic and markers effects together led to a better global fit and predictive ability to breeding models, especially when penalized methods are applied to markers effects.

### Conflict of interest statement

No conflict declared.

### Acknowledgments

We are grateful to the Fundação de Amparo a Pesquisa do Estado de São Paulo (FAPESP), Merial/Igenity and Conselho Nacional de apoio a Pesquisa (CNPq) for the financial support, to Agro-Pecuária CFM for data set and the Instituto de Investigación y Tecnología Agroalimentarias de Cataluña (IRTA) as the host institution for its full backing while preparing the research and the manuscript.

### References

- Ayres, D.R., Souza, F.R.P., Mercadante, M.E.Z., Fonseca, L.F.S., Tonhati, H., et al., 2010. Evaluation of TFAM and FABP4 gene polymorphisms in three lines of Nellore cattle selected for growth. *Genet. Mol. Res.* 9 (4), 2050–2059.
- Blasco, A., Piles, M., Varona, L., 2003. A Bayesian analysis of the effect of selection for growth rate on growth curves in rabbits. *Genet. Sel. Evol.* 35, 21–41.
- Boldman, K.G., Kriese, L.A., Van Vleck, L.D., Kachman, S.D., 1995. A Manual for Use of MTDPREML: A Set of Programs to Obtain Estimates of Variances and Covariances. United States Department of Agriculture-Agricultural Research Service, Nebraska (pp. 115).
- De los Campos, G., Naya, H., Gianola, D., Crossa, J., Legarra, A., et al., 2009. Predicting quantitative traits with regression models for dense molecular markers and pedigree. *Genetics* 182, 375–385.
- Fernando, R.L., Grossman, M., 1989. Marker assisted selection using best linear unbiased prediction. *Gen. Sel. Evol.* 21, 467–477.
- Ferraz, J.B.S., Pinto, L.F.B., Meirelles, F.V., Eler, J.P., Rezende, F.M., et al., 2009. Association of single nucleotide polymorphisms with carcass traits in Nellore cattle. *Genet. Mol. Res.* 8 (4), 1360–1366.
- Geldermann, H., 1975. Investigation on inheritance of quantitative characters in animals by gene markers. 1. Methods. *Theor. Appl. Genet.* 46, 319–333.
- Geweke, J., 1992. Evaluating the accuracy of sampling-based approaches to the calculation of posterior moments (with discussion). In: Bernardo, J.M., Berger, J.O., Dawid, A.P., Smith, A.F.M. (Eds.), *Bayesian Statistics*, Oxford University Press, Oxford, GB, pp. 169–193.
- Geyer, C.J., 1992. Practical Markov chain Monte Carlo. *Stat. Sci.* 7, 473–511.
- Gianola, D., Perez-Enciso, M., Toro, M.A., 2003. On marker-assisted prediction of genetic value: beyond the ridge. *Genetics* 163, 347–365.
- Gianola, D., Fernando, R.L., Stella, A., 2006. Genomic-assisted prediction of genetic value with semiparametric procedures. *Genetics* 173, 1761–1776.
- Habier, D., Fernando, R.L., Garrick, D.J., 2011. Extension of the Bayesian alphabet for genomic selection. *BMC Bioinformatics* 12, 186.
- González-Recio, O., Gianola, D., Long, N., Weigel, K.A., Rosa, G.J.M., et al., 2008. Nonparametric methods for incorporating genomic information into genetic evaluations: an application to mortality in broilers. *Genetics* 178, 2305–2313.
- Henderson, C.R., 1984. In: *Applications of Linear Models in Animal Breeding* University of Guelph, Ontario (pp. 462).
- Horimoto, A.R.V.R., Ferraz, J.B.S., Balieiro, J.C.C., Eler, J.P., 2007. Phenotypic and genetic correlations for body structure scores (frame) with productive traits and index for CEIP classification in Nellore beef cattle. *Genet. Mol. Res.* 6 (1), 188–196.
- Ibañez-Escriche, N., Gonzalez-Recio, O., 2011. Review. Promises, pitfalls and challenges of genomic selection in breeding programs. *Span. J. Agric. Res.* 9 (2), 404–413.
- Kendall, M.G., 1947. In: *The Advanced Theory of Statistics* Charles Griffin & Company Limited, London (pp. 457).
- Lande, R., Thompson, R., 1990. Efficiency of marker-assisted selection in the improvement of quantitative traits. *Genetics* 124, 743–756.
- Legarra, A., Aguilar, I., Misztal, I., 2009. A relationship matrix including full pedigree and genomic information. *J. Dairy Sci.* 92, 4656–4663.
- Legarra, A., Varona, L., López de Maturana, E., 2008. TM: Threshold Model. <<http://snp.toulouse.inra.fr/~alegarra>>.
- Meuwissen, T.H.E., Hayes, B.J., Goddard, M.E., 2001. Prediction of total genetic value using genome-wide dense marker maps. *Genetics* 157, 1819–1829.
- Pinto, L.F.B., Ferraz, J.B.S., Meirelles, F.V., Eler, J.P., Rezende, F.M., et al., 2010. Association of SNPs on CAPN1 and CAST genes with tenderness in Nellore cattle. *Genet. Mol. Res.* 9 (3), 1431–1442.
- Raftery, A.E., Lewis, S.M., 1992. How many iterations in the Gibbs Sampler?. In: Bernardo, J.M., Berger, J.O., Dawid, A.P., Smith, A.F.M. (Eds.), *Bayesian Statistics*, vol. IV. , Oxford University Press, New York, USA, pp. 763–774.
- Soller, M., 1978. The use of loci associated with quantitative effects in dairy cattle improvement. *Anim. Prod.* 27, 133–139.
- Spiegelhalter, D.J., Best, N.G., Carlin, B.P., Van derLinde, A., 2002. Bayesian measures of model complexity and fit (with discussion). *J. R. Stat. Soc. Ser. B* 64, 583–639.
- Van Melis, M.H., Oliveira, H.N., Eler, J.P., Ferraz, J.B.S., Casellas, J., Varona, L., 2010. Additive genetic relationship of longevity with fertility and production traits in Nellore cattle based on bivariate models. *Genet. Mol. Res.* 9 (1), 176–187.
- VanRaden, P.M., 2008. Efficient methods to compute genomic predictions. *J. Dairy Sci.* 91, 4414–4423.
- Whittaker, J.C., Thompson, R., Denham, M.C., 2000. Marker-assisted selection using ridge regression. *Genet. Res.* 75, 249–252.