



# Joint genome-wide prediction in several populations accounting for randomness of genotypes: A hierarchical Bayes approach. I: Multivariate Gaussian priors for marker effects and derivation of the joint probability mass function of genotypes

Carlos Alberto Martínez<sup>a,b,\*</sup>, Kshitij Khare<sup>b</sup>, Arunava Banerjee<sup>c</sup>, Mauricio A. Elzo<sup>a</sup>

<sup>a</sup> Department of Animal Sciences, University of Florida, Gainesville, FL, USA

<sup>b</sup> Department of Statistics, University of Florida, Gainesville, FL, USA

<sup>c</sup> Department of Computer and Information Science and Engineering, University of Florida, Gainesville, FL, USA

## ARTICLE INFO

### Keywords:

Across population genome-enabled prediction  
Bayesian modeling  
Heterogeneous allelic frequencies  
Distribution of genotypes

## ABSTRACT

It is important to consider heterogeneity of marker effects and allelic frequencies in across population genome-wide prediction studies. Moreover, all regression models used in genome-wide prediction overlook randomness of genotypes. In this study, a family of hierarchical Bayesian models to perform across population genome-wide prediction modeling genotypes as random variables and allowing population-specific effects for each marker was developed. Models shared a common structure and differed in the priors used and the assumption about residual variances (homogeneous or heterogeneous). Randomness of genotypes was accounted for by deriving the joint probability mass function of marker genotypes conditional on allelic frequencies and pedigree information. As a consequence, these models incorporated kinship and genotypic information that not only permitted to account for heterogeneity of allelic frequencies, but also to include individuals with missing genotypes at some or all loci without the need for previous imputation. This was possible because the non-observed fraction of the design matrix was treated as an unknown model parameter. For each model, a simpler version ignoring population structure, but still accounting for randomness of genotypes was proposed. Implementation of these models and computation of some criteria for model comparison were illustrated using two simulated datasets. Theoretical and computational issues along with possible applications, extensions and refinements were discussed. Some features of the models developed in this study make them promising for genome-wide prediction, the use of information contained in the probability distribution of genotypes is perhaps the most appealing. Further studies to assess the performance of the models proposed here and also to compare them with conventional models used in genome-wide prediction are needed.

## 1. Introduction

The use of molecular markers located across the whole genome for prediction of breeding values (Meuwissen et al., 2001) and phenotypes (Goddard and Hayes, 2007; Gianola et al., 2009) has proven to be a useful tool in animals (Hayes et al., 2009), humans (Guttmacher et al., 2002; de los Campos et al., 2010) and plants (Bernardo and Yu, 2007; Desta and Ortiz, 2014). This success has given rise to a tremendous amount of research in the area of statistical genomics in order to obtain better genome-wide predictions (Goddard and Hayes, 2007; Gianola, 2013; Hill, 2014; Gianola and Rosa, 2015).

Most of the methods have been developed for prediction in a single population. Across population studies usually use predictions obtained

from individual populations or pool data to perform a single analysis (de Roos et al., 2009). On one hand, pooling data and performing a single analysis may increase the accuracy of genome-wide prediction because the number of records has an important impact on it (Meuwissen et al., 2001; Goddard, 2009; Zhong et al., 2009). On the other hand, it may decrease accuracy when the effects of QTL controlling the trait are not the same across populations (de Roos et al., 2009; van den Berg et al., 2015; Wientjes et al., 2015).

Analyzing data from Holstein cattle performing in different European countries, Lund et al. (2011) reported that pooling data and carrying out a single analysis increased the accuracy of genomic predictions. With simulated data, de Roos et al. (2009) found that pooling data was beneficial when populations had diverged by few

\* Corresponding author at: Department of Animal Sciences, University of Florida, Gainesville, FL 32611, USA.  
E-mail address: [carlosmn@ufl.edu](mailto:carlosmn@ufl.edu) (C.A. Martínez).

generations, marker density was high and heritability was low, but for more distant populations and less dense marker panels they found a small decrease in accuracy. Using simulated data, Wientjes et al. (2015) studied the effect of differences in QTL allele substitution effects across populations on the accuracy of genome-wide prediction. They found that when allele substitution effects changed across populations, the accuracies decreased in proportion to the genetic correlation between populations. Using the same dataset, van den Berg et al. (2015) looked for across population genomic prediction scenarios under which Bayesian variable selection models had a better performance than genomic BLUP (GBLUP). They concluded that Bayesian variable selection models outperform GBLUP when the number of QTL is small as in single population analyses, but the difference in accuracy is larger in the across population case.

None of these studies allowed marker effects to differ from one population to another. However, de Roos et al. (2009) highlighted the need for alternative methods that allow population-specific estimation of allele substitution effects in across population genome wide prediction. Chen et al. (2014) proposed a Bayesian model with different SNP effects for each population that permits sharing information across populations through a common set of latent variables indicating whether a given marker is associated with a QTL or not. They did not model covariance matrices of marker effects explicitly. With real and simulated data they found that this model increased the accuracy of across population genome-wide prediction, especially when the number of QTL was small and correlations among QTL effects from different populations were high. Recently, Bayesian models that account for genetic heterogeneity have been proposed. Multivariate models considering correlated population specific marker effects were developed by Lehermeier et al. (2015) while de los Campos et al. (2015a) proposed a model with main marker effects and interactions. Using real data from three plant populations, Lehermeier et al. (2015) found cases in which the strategy of pooling data and ignoring structure performed better and others where the multivariate models yielded better predictive performance. For example, in highly differentiated populations within group and multivariate analyses performed better. Using real datasets from pigs and wheat, de los Campos et al. (2015a) found modest superiority of the interaction model relative to the model using pooled data and the model that analyzed each subpopulation separately. Similar studies have implemented multivariate models in multibreed dairy cattle populations (Karoui et al., 2012; Olson et al., 2012; Makgahlela et al., 2013). Huang et al. (2014) used non-linear models to perform genome wide prediction in layer hens when the reference population was comprised by individuals from several breeds or lines and compared them with a multiple-trait GBLUP model. They found that the various models used had a similar predictive performance.

If several populations are to be evaluated simultaneously, the possible existence of genotype by environment interaction, lack of persistence of linkage phase and variation in allelic frequencies across populations indicate the need for an analysis that accounts for the fact that combining them creates a structured complete population. It has been reported that population structure may act as an effect modifier (de los Campos et al., 2015a). Furthermore, it has to be considered that not only the allele substitution effects of a particular locus in different populations may be correlated, but also its frequencies in each population (e.g., due to gene flow).

Another feature that has been overlooked in the random linear regression models used in genome-wide prediction is the randomness of the matrix containing a one to one mapping from the set of genotypes to a subset of the integers, namely the design matrix. This matrix is treated as fixed in genome-wide prediction models, while in classical quantitative genetics theory it is treated as random (Falconer and Mackay, 1996; Lynch and Walsh, 1998). Besides being in agreement with the classical theory, taking into account the randomness of this matrix, that is, the randomness of genotypes, permits the estima-

tion of allelic frequencies because when treated as an observable discrete random matrix, its probability mass function (pmf) depends on the allelic frequencies. Thus, under a Bayesian setting, allelic frequencies are treated as random because these are unknown parameters. Further, the works of Wright (1930, 1937) provide additional support to treat allelic frequencies as random variables making Bayesian inference even more attractive.

Thus, the objective of this study was to propose hierarchical Bayesian models to carry out simultaneous genome-wide prediction in several populations accounting for randomness of marker genotypes, heterogeneity and correlation of allelic frequencies across populations, and population-specific allelic substitution effects.

## 2. Methods

### 2.1. The models

Hereinafter the complete population or simply the population is defined as the set of individuals with phenotypes considered in the study. Suppose that there exists some criterion (e.g., environment, race, breed, line, etc.) to split this population into  $S$  subpopulations. To make the problem more tractable, some simplifying assumptions are made. The first one is linkage equilibrium. The second one is Hardy-Weinberg equilibrium. The third one is that starting from the oldest individuals with phenotypes, the pedigree is fully known. Lastly, mutations are ignored.

The basic linear model used to describe the relationship between response variables and marker allele substitution effects is  $y = Wg + e$ , where  $y$  is a vector containing dependent variables (e.g., records corrected for non-genetic factors),  $W$  is an observable random matrix containing a one to one mapping from individual marker genotypes to a subset of the integers to be defined later,  $g$  is an unknown random vector of marker allelic substitution effects for every population and  $e$  is a random vector of residuals. A more detailed notation is the following. If records are sorted by subpopulation as well as the columns of  $W$  and the elements of  $g$ , then for every  $l = 1, 2, \dots, S$ ,  $y_l = W_l g_l + e_l$ , with dimensions:  $(y_l)_{n_l \times 1}$ ,  $(W_l)_{n_l \times m}$ ,  $(g_l)_{m \times 1}$  and  $(e_l)_{n_l \times 1}$  where  $n_l$  is the sample size of subpopulation  $l$ , and  $m$  is the number of marker loci. Thus, the total sample size is  $n = \sum_{l=1}^S n_l$ .

The scenario where only a part of matrix  $W$  is observed because some individuals are not genotyped or individuals are genotyped for different numbers of marker loci is also considered. This is done by treating this non-observed part of  $W$  as a parameter in the model as it will be explained later.

The case of diploid individuals and biallelic marker loci is considered. The effect of every marker locus is defined as the regression of records on a function of the number of copies of the reference allele and in quantitative genetics it corresponds to the allele substitution effect (Falconer and Mackay, 1996; Lynch and Walsh, 1998). The number of copies can be “centered” at zero giving the following codification. Let  $A$  and  $B$  be the marker alleles at each locus and let  $B$  be the reference allele. Then:

$$W_l = \{w_{ij}^l\}_{n_l \times m} = \begin{cases} 1, & \text{if genotype} = BB \\ 0, & \text{if genotype} = AB \\ -1, & \text{if genotype} = AA \end{cases}$$

Different versions of the hierarchy that represents the stochastic component of each model were considered. Models vary according to the assumptions on the variance of residuals and the priors posed over the marker effects. The most parsimonious model is the one considering homoscedastic residuals and homogeneous marker effect covariance matrices. The hierarchical Bayesian model assuming homoscedastic residuals and multivariate Gaussian priors for marker effects has the following structure:

$$y|W, g, \sigma^2 \sim MVN(Wg, \sigma^2 I)$$

$$W|p_1^*, p_2^*, \dots, p_m^* \sim \pi(p_1^*, p_2^*, \dots, p_m^*)$$

$$p_j^* \stackrel{iid}{\sim} \pi(p^*), \quad j=1, 2, \dots, m$$

$$\sigma^2 \sim \text{InverseGamma}\left(\frac{\tau^2}{2}, \frac{v}{2}\right) := IG\left(\frac{\tau^2}{2}, \frac{v}{2}\right)$$

$$g|G \sim \text{MVN}(0, G), \quad G = \text{BlockDiag}\{G_j\}_{j=1}^m$$

$$G_j \stackrel{iid}{\sim} \text{InverseWishart}(a, \Sigma) := IW(a, \Sigma)$$

$$G_j = \begin{bmatrix} \sigma_{j1}^2 & \sigma_{j1,2} & \cdots & \sigma_{j1,S} \\ & \sigma_{j2}^2 & \cdots & \sigma_{j2,S} \\ & & \ddots & \vdots \\ \text{sym} & & & \sigma_{jS}^2 \end{bmatrix}$$

where  $\sigma^2$  is the residual variance,  $\sigma_{jl}^2$  is the variance of the effect of the  $j^{\text{th}}$  marker in the  $l^{\text{th}}$  subpopulation,  $\sigma_{jll'}$  is the covariance between effects of marker  $j$  in subpopulations  $l$  and  $l'$ ,  $p_j^*$  is a parameter associated with allelic frequencies of the  $j^{\text{th}}$  marker in each subpopulation and  $\pi(p^*)$  is its density. Details on these parameters and their probability density function (pdf) are given later.

In the case of heterogeneous residual variances across subpopulations, residual variances  $\sigma_1^2, \dots, \sigma_S^2$  are given independent  $IG\left(\frac{\tau^2}{2}, \frac{v}{2}\right)$  priors and then:  $y|W, g, R \sim \text{MVN}(Wg, V)$ ,  $R = (\sigma_{e1}^2, \dots, \sigma_{eS}^2)$  and  $V = \text{BlockDiag}\{\sigma_{ei}^2 I_{n_i}\}_{i=1}^S$ . Hill (1984) found that in the presence of heterogeneous environmental variances, across population analyses assuming homogenous residuals variances yielded an excess of individuals selected from populations with higher environmental variances. This is why heterogeneity of residual variances across subpopulations was considered in this study.

The general framework assumes that in each subpopulation there is a fraction of genotyped individuals and a fraction of non-genotyped or partially genotyped individuals. Let  $W^o$  and  $W^N$  denote the observed (data) and non-observed (an unknown parameter) parts of  $W$ . Let  $P^* = (p_1^*, p_2^*, \dots, p_m^*)$ ; therefore,  $\pi(W|P^*) = \pi(W^o, W^N|P^*)$  can be expressed as:  $f(W^o|W^N, P^*)\pi(W^N|P^*)$ . Thus, the full likelihood has the form:

$$f(y, W^o|W^N, g, R, P^*) = f(y|W^o, W^N, g, R, P^*)f(W^o|W^N, g, R, P^*)$$

$$= f(y|W, g, R)f(W^o|W^N, P^*).$$

Henceforth,  $f(y|W, g, R)$  will be referred to as the  $y$  component of the likelihood and  $f(W^o|W^N, P^*)$  will be referred to as the  $W$  component.

The simplest case for the covariance matrix of marker effects is  $G = I \otimes G^0$ . Under this setting the assumption is that the covariance structure is the same for all markers. This is statistically convenient due to the fact that the number of covariance parameters is reduced. Further, in analysis considering a single population, it has been found that specifying a different variance for each marker does not allow too much Bayesian learning about marker effect variances (Gianola et al., 2009). Here, models assigning the same covariance matrix to the effects of all marker loci and models considering a different covariance matrix for the effects of each marker locus were considered and these models were referred to as homogeneous marker effect covariance matrix models and heterogeneous marker effect covariance matrix models. Let  $\mathcal{P}_S^+$  denote the space of symmetric positive definite matrices of dimension  $S \times S$ . Then, the marginal prior distribution of  $g$  is:

$$\pi(g) = \int_{\mathcal{P}_S^+} \pi(g|G^0)\pi(G^0)dG^0 \propto \frac{1}{\left|\Sigma + \sum_{j=1}^m g_j g_j'\right|^{\left(\frac{a+m}{2}\right)}}.$$

For details, see Appendix A. Similarly, for the heterogeneous marker effect covariance matrix model it can be shown (appendix A) that:  $\pi(g) \propto \frac{1}{\prod_{j=1}^m \left(1 + \frac{1}{a+1-S} g_j' \Sigma^{-1} g_j\right)^{\left(\frac{a+1}{2}\right)}}$ , which is the product of  $m$  multi-

variate  $t$  distributions with scale matrix  $\Sigma_* = \frac{1}{a+1-S} \Sigma$  and degrees of freedom  $a+1-S$ ; therefore, under this prior, marker effects are marginally independent and identically distributed. At this point, the following remark can be made.

**Remark 1.** Under the assumption of homogeneous marker effect covariance matrices, *a priori* the marker effects are marginally dependent. This happens because when integrating with respect to the common covariance matrix  $G^0$ , the term  $\sum_{j=1}^m g_j g_j'$  and the hyper-hyperparameter  $\Sigma$  are factored, resulting in a function that cannot be written as the product of  $m$  functions, each one depending on a different  $g_j$ . Moreover, the joint prior density is not standard.

To take into account the belief that allelic frequencies of the same marker vary across subpopulations and may be correlated, the prior  $\pi(p^*)$  is built based on a Dirichlet distribution. To do that, the allelic frequency of the reference allele in marker locus  $j$  in subpopulation  $l$  has to be expressed on a complete population basis, that is,  $p_{lj}$  is expressing the frequency of the reference allele in locus  $j$  in subpopulation  $l$  relative not to subpopulation  $l$ , but to the complete population. Thus, the frequencies of the two alleles at a given marker locus and a given subpopulation do not add to one, but to some sort of relative frequency of that subpopulation in that locus denoted as  $r_{lj}$ . Let  $r = (r_1, \dots, r_S)$ ,  $r_l = (r_{l1}, \dots, r_{lm})$ ,  $l = 1, 2, \dots, S$ . With this parameterization  $\sum_{j=1}^m p_{lj} \leq 1$ ,  $\forall j = 1, 2, \dots, m$ , with equality if and only if the reference allele is fixed in all subpopulations. Conversely, allelic frequencies expressed on a subpopulation basis satisfy the constraint that the sum of the frequencies of the two alleles at each marker locus equals one within each subpopulation. Let  $q_{lj} = 1, 2, \dots, m$ ,  $l = 1, 2, \dots, S$ , be the frequencies of the non-reference alleles expressed on a complete population basis, then  $p_{lj} + q_{lj} = r_{lj}$ . The two parameterizations of allelic frequencies are related by the one to one mapping  $p_{lj}^* = p_{lj}/r_{lj}$ .

Consider the case when  $r$  is known and  $r_{11} = \dots = r_{lm} = r_l \forall l$ . Then, elements of vector  $r = (r_1, \dots, r_S)$  can be seen as subpopulation weights, that is, they are related to subpopulation sizes. By  $r$  being known, it is meant that it is either actually known or it is specified following some assumption. A pragmatic decision would be to assign equal subpopulation weights, an assumption that was also made in other studies (e.g., Gianola et al., 2010). Once  $r$  has been specified, there is an extra restriction over each  $p_j = (p_{1j}, \dots, p_{Sj})$ . For  $l = 1, 2, \dots, S$  the following condition must be satisfied:  $p_{lj} \leq r_l$ . Therefore, the support of the distribution of  $p_j$  given  $r$  is  $\Omega_j^r = \{p_j \in \mathbb{R}^S | 0 < p_{lj} \leq r_l \forall l, \sum_{l=1}^S r_l = 1\}$ . Notice that the condition  $\sum_{l=1}^S r_l = 1$  implies that vectors in  $\Omega_j^r$  satisfy  $\sum_{l=1}^S p_{lj} \leq 1$ . Thus, under this approach the prior used for each  $p_j$  is one corresponding to a scaled Dirichlet random vector. If  $\beta = (\beta_1, \dots, \beta_S) \sim \text{Dirichlet}(\alpha)$ ,  $\alpha \in \mathbb{R}^{S+1}$ , then the prior assigned to  $p_j$  is the distribution of vector  $(\beta_1 r_1, \dots, \beta_S r_S)$  which clearly pertains to  $\Omega_j^r$ . Then, the pdf  $\pi(p_j|r)$  is derived using standard results from the theory of distributions of transformations of random variables (Casella and Berger, 2002). This derivation is simplified by the fact that the transformation is linear and therefore the Jacobian is constant. It follows that:

$$\pi(p_j|r) \propto \prod_{l=1}^S \left\{ \left( \frac{p_{lj}}{r_l} \right)^{\alpha_l - 1} \right\} p_{(S+1)j}^{\alpha_{S+1} - 1}, \text{ where } p_{(S+1)j} = 1 - \sum_{l=1}^S \frac{p_{lj}}{r_l}.$$

The second approach is to assume that  $r$  is unknown. The density  $\pi(p|r)$  could be used and a Dirichlet distribution could be assigned to each  $r_j$  adding one more level to the hierarchy. However, using  $p_{lj} + q_{lj} = r_{lj}$  and properties of the Dirichlet distribution, the following strategy allows assigning a prior to allelic frequencies and the weights  $r$  without putting an extra level in the hierarchy. To this end it is

assumed that  $r_{lj}$  varies for each  $j$  and each  $l$ . A *Dirichlet*  $((\alpha_p, \alpha_q))$  prior is posed over  $(p_j, q_j)$ , where  $q_j$  is the analog of  $p_j$  for the non-reference allele at each locus and  $\alpha_p = (\alpha_{1p}, \dots, \alpha_{Sp})$ ,  $\alpha_q = (\alpha_{1q}, \dots, \alpha_{Sq})$ . Consequently, by properties of the Dirichlet distribution it follows that  $r_j \sim \text{Dirichlet}((\alpha_{1p} + \alpha_{1q}, \dots, \alpha_{Sp} + \alpha_{Sq}))$ .

### 2.1.1. Deriving the joint pmf of marker genotypes conditional on allelic frequencies

Given the kinship structure of a population (i.e., the pedigree) one can find several generations comprised of genotyped, partially genotyped and non-genotyped individuals. Therefore, the approach is to derive the pmf of the complete matrix  $W$ , i.e., the joint pmf of individuals with phenotypic records. Under this setting,  $m$  is the total number of marker loci to be included in the analysis (it usually corresponds to the size of the densest marker panel used in the population).

Across columns, that is, across marker loci, the problem is simplified by assuming linkage equilibrium, which implies independence of genotypes at different loci. Therefore, for an arbitrary subpopulation, the joint density of its column vectors is simply the product of their marginal pmf. When considering all subpopulations, the same assumption implies that marker genotypes at different loci are independent. The following derivations hold for any of the previously discussed approaches to model allelic frequencies distributions. Under the assumption of Hardy-Weinberg equilibrium it follows that marginally:

$$w_{ij}^l p_{lj}^* \sim \begin{cases} 1, & \text{with probability } p_{lj}^{*2} \\ 0, & \text{with probability } 2p_{lj}^*(1 - p_{lj}^*) \\ -1, & \text{with probability } (1 - p_{lj}^*)^2 \end{cases}$$

Recall that  $p_{lj}^* = p_{lj}/r_{lj}$ . Notice that  $p_{lj}^*$  is used instead of  $p_{lj}$  because it allows defining a proper pmf in the sense that the sum of the probabilities of the three possible values of  $w_{ij}^l$  equals one (which does not happen when using  $p_{lj}$ ). The pmf  $\pi(w_{ij}^l p_{lj}^*)$  can be also written as:

$$\pi(w_{ij}^l p_{lj}^*) = (p_{lj}^{*2})^{I_{ij}} (2p_{lj}^*(1 - p_{lj}^*))^{I_{0i}} ((1 - p_{lj}^*)^2)^{I_{-i}},$$

where  $I_{ij}$  is the indicator variable of the mutually exclusive events  $w_{ij}^l = z$ ,  $z \in \{-1, 0, 1\}$ . By the linkage equilibrium assumption it follows that for individual  $i$  in population  $l$ :  $\pi(w_i^l p_j^*) = \prod_{j=1}^m \pi(w_{ij}^l p_{lj}^*)$ .

The rows of matrix  $W$  represent individuals with records. Because of the kinship between them, the genotype of a given individual is not independent of the genotype of their relatives. Furthermore, this non-independence has to be considered across subpopulations (e.g., half or full sibs may pertain to different subpopulations). This approach is based on the pedigree of the complete population. The “base” animals or “founders” can be pragmatically defined as the oldest individuals with phenotypic records and those individuals with phenotypes and unknown parents. To facilitate computations, it is assumed that these individuals are unrelated. Hereinafter this set is referred to as the base population, and individuals in this set are referred to as founders or base individuals. The remaining individuals in the population are referred to as non-founders. This pmf could be derived ignoring pedigree information which is equivalent to mutual independence of the rows of  $W$ , then  $\pi(W|P^*) = \prod_{j=1}^m \prod_{l=1}^S \prod_{i=1}^{n_l} \pi(w_{ij}^l p_{lj}^*)$ . However, this would ignore information contained in the pedigree and would unnecessarily make the parametric space of  $W^N$  larger, which does not seem to be the best way to proceed.

The ordering of individuals is arbitrary, but a convenient way to do it here is according to the pedigree in such a way that the founders are given the first indices. For marker locus  $j$  in population  $l$  the target is to find:

$$\pi(w_j^l p_j^*) = \pi(w_{1j}^l, w_{2j}^l, \dots, w_{n_{lj}}^l p_{lj}^*) = P(w_{1j}^l = \omega_1, w_{2j}^l = \omega_2, \dots, w_{n_{lj}}^l = \omega_{n_l} | p_j^*)$$

with  $\omega_i \in \{-1, 0, 1\}$ ,  $1 \leq i \leq n_l$ . This joint pmf can be written as:

$$\begin{aligned} \pi(w_j^l p_j^*) &= \pi(w_{n_{lj}}^l | w_{1j}^l, \dots, w_{(n_l-1)j}^l, p_j^*) \pi(w_{(n_l-1)j}^l, \dots, w_{(n_l-2)j}^l p_j^*) \\ &= \pi(w_{n_{lj}}^l | w_{1j}^l, \dots, w_{(n_l-1)j}^l, p_j^*) \pi(w_{(n_l-1)j}^l | w_{1j}^l, \dots, w_{(n_l-2)j}^l, p_j^*) \\ &\quad \times \pi(w_{1j}^l, \dots, w_{(n_l-2)j}^l p_j^*) \\ &= \pi(w_{n_{lj}}^l | w_{1j}^l, \dots, w_{(n_l-1)j}^l, p_j^*) \dots \pi(w_{1j}^l p_j^*) \\ &= \prod_{i=0}^{n_l-2} \{ \pi(w_{(n_l-i)j}^l | w_{1j}^l, \dots, w_{(n_l-i-1)j}^l, p_j^*) \} \pi(w_{1j}^l p_j^*). \end{aligned}$$

When considering all the  $m$  marker loci we have:

$$\pi(W|P^*) = \prod_{i=0}^{n_l-2} \{ \pi(w_{(n_l-i)j}^l | w_{1j}^l, \dots, w_{(n_l-i-1)j}^l, p_j^*) \} \pi(w_{1j}^l p_j^*),$$

where each one of the pmf  $\pi(w_{(n_l-i)j}^l | w_{1j}^l, \dots, w_{(n_l-i-1)j}^l, p_j^*)$  is the product:  $\prod_{j=1}^m \pi(w_{(n_l-i)j}^l | w_{1j}^l, \dots, w_{(n_l-i-1)j}^l, p_j^*)$ ,  $0 \leq i \leq n_l-2$  and  $\pi(w_{1j}^l p_j^*) = \prod_{j=1}^m \pi(w_{1j}^l p_{lj}^*)$ .

Now, a conditional independence argument is used to simplify  $\pi(W|P^*)$ . Given the genotypes of the parents of individual  $i$ , its genotype is independent of the genotype of collateral relatives and other ancestors. It is possible that the parents of individual  $i$  in population  $l$  pertain to subpopulations  $l^*$  and  $l'$ . Thus, at this point the complete population is considered. In addition, notice that given the parental genotypes, the genotype of an individual does not depend on the allelic frequencies because this conditional pmf is determined using basic segregation rules (see Appendix A). From these arguments it follows that for individual  $i$ ,  $\pi(w_i | w_1, \dots, w_{i-1}, P^*) = \pi(w_i | w_S, w_{D_i})$ , where  $w_S$  and  $w_{D_i}$  are the genotypes of the parents of individual  $i$ . The pmf of non-founder genotypes at marker locus  $j$  conditioned on their parental genotypes is presented in Appendix A. Therefore,  $\pi(W|P^*)$  can be written as  $\pi(W|P^*) = \pi(W_{NF}|W_F) \pi(W_F|P^*)$  where  $W_F$  is the submatrix of  $W$  formed by considering the rows corresponding to founders and  $W_{NF}$  is the submatrix of  $W$  comprised of the rows corresponding to non-founders. Let  $f$  be the total number of founders. Under the assumption that these individuals are unrelated, the pmf of their genotypes given allelic frequencies is:

$$\begin{aligned} \pi(W_F|P^*) &= \prod_{i=1}^f \pi(W_i|P^*) = \prod_{j=1}^m \prod_{i=1}^f \pi(w_{ij} | P^*) = \prod_{j=1}^m \prod_{l=1}^S \prod_{i=1}^{f_l} \pi(w_{ij}^l p_{lj}^*) \\ &= \prod_{j=1}^m \prod_{l=1}^S \prod_{i=1}^{f_l} (p_{lj}^{*2})^{I_{ij}} (2p_{lj}^*(1 - p_{lj}^*))^{I_{0i}} ((1 - p_{lj}^*)^2)^{I_{-i}} \\ &= \prod_{j=1}^m \prod_{l=1}^S (p_{lj}^{*2})^{n_l^{BBj}} (2p_{lj}^*(1 - p_{lj}^*))^{n_l^{ABj}} ((1 - p_{lj}^*)^2)^{n_l^{AAj}} \\ &= \prod_{j=1}^m \prod_{l=1}^S 2^{n_l^{ABj}} p_{lj}^{*2n_l^{BBj} + n_l^{ABj}} (1 - p_{lj}^*)^{2n_l^{AAj} + n_l^{ABj}} = 2^{n^H} \prod_{j=1}^m \prod_{l=1}^S p_{lj}^{*2n_l^{BBj}} (1 - p_{lj}^*)^{n_l^{ABj}}, \end{aligned}$$

replacing

$$p_{lj}^* = p_{lj}/r_{lj} \quad \forall \quad l=1, 2, \dots, S, \forall j=1, 2, \dots, m:$$

$$\pi(W_F|P, r) = 2^{n^H} \prod_{j=1}^m \prod_{l=1}^S \frac{1}{r_{lj}^{2n_l^{BBj}}} p_{lj}^{2n_l^{BBj}} (r_{lj} - p_{lj})^{n_l^{ABj}} \text{ where } f_l \text{ is the number of}$$

founders in the  $l^{\text{th}}$  subpopulation; thus,  $f = \sum_{l=1}^S f_l$ ,  $n_l^{BBj}$ ,  $n_l^{ABj}$  and  $n_l^{AAj}$  are the counts of founders with genotypes BB, AB and AA at marker locus  $j$  in subpopulation  $l$  respectively,  $n_l^{BBj} = 2n_l^{BBj} + n_l^{ABj}$  is the total count of B alleles at marker locus  $j$  in founders from subpopulation  $l$ ,  $n_l^{Aj} = 2n_l^{AAj} + n_l^{ABj}$  is the total count of A alleles at marker locus  $j$  in founders from subpopulation  $l$  and  $n^H = \sum_{j=1}^m \sum_{l=1}^S n_l^{ABj}$  is the total number of heterozygous loci in the base population. In terms of the random variables  $w_{ij}^l$ ,  $n_l^{BBj}$ ,  $n_l^{ABj}$  and  $n_l^{AAj}$  can be written as:  $n_l^{BBj} = \sum_{i=1}^{f_l} I_{li}$ ,  $n_l^{AAj} = \sum_{i=1}^{f_l} I_{-li}$ ,  $n_l^{ABj} = f_l - (n_l^{BBj} + n_l^{AAj}) = f_l - \sum_{i=1}^{f_l} (w_{ij}^l)^2$ . For non-founders:



$$\pi(W_{NF}|W_F) = \prod_{j=1}^m \prod_{i'=f+1}^n \pi(w_{i'j}^l | w_{S_{i'j}}, w_{D_{i'j}}) = \prod_{j=1}^m \prod_{l=1}^S \prod_{i'=f_l+1}^{n_l} \pi(w_{i'j}^l | w_{S_{i'j}}, w_{D_{i'j}})$$

where  $w_{S_{i'j}}^l$  and  $w_{D_{i'j}}^l$  are the genotypes for marker  $j$  of the parents of individual  $i'$  from subpopulation  $l$ . Hence:

$$\begin{aligned} \pi(W|P^*) &= \prod_{j=1}^m \prod_{l=1}^S \prod_{i=1}^{f_l} \pi(w_{ij}^l | p_{ij}^{*l}) \times \prod_{j=1}^m \prod_{l=1}^S \prod_{i'=f_l+1}^{n_l} \pi(w_{i'j}^l | w_{S_{i'j}}, w_{D_{i'j}}) \\ &= \prod_{j=1}^m \prod_{l=1}^S \prod_{i=1}^{f_l} \left\{ \pi(w_{ij}^l | p_{ij}^{*l}) \times \prod_{i'=f_l+1}^{n_l} \pi(w_{i'j}^l | w_{S_{i'j}}, w_{D_{i'j}}) \right\} \\ &= 2^{nH} \prod_{j=1}^m \prod_{l=1}^S \left\{ p_{ij}^{*l} p_{ij}^{*l} (1-p_{ij}^{*l})^{n_l} \prod_{i'=f_l+1}^{n_l} \pi(w_{i'j}^l | w_{S_{i'j}}, w_{D_{i'j}}) \right\} \\ &\Rightarrow \pi(W|P, r) = 2^{nH} \prod_{j=1}^m \prod_{l=1}^S \left\{ \frac{1}{r_{ij}^{2f_l}} p_{ij}^{*l} p_{ij}^{*l} (r_{ij} - p_{ij}^{*l})^{n_l} \prod_{i'=f_l+1}^{n_l} \pi(w_{i'j}^l | w_{S_{i'j}}, w_{D_{i'j}}) \right\}. \end{aligned}$$

**Remark 2.** Under the assumptions presented at the beginning of this section, given base genotypes, the process defining the inheritance of alleles is completely determined by the pedigree information. The pedigree allows tracing the set of possible values that genotypes can take from a given individual back to the base population. It implies that allelic frequencies have to be known only in the base population because the distribution of genotypes in the set of non-founders is completely determined by the pedigree. Stated another way, given the pedigree, only the founder genotypes carry information about allelic frequencies.

The next step is to formally define the support (set of values of  $W$  with non-null probability) of the pmf  $\pi(W|P^*)$  and its cardinality (i.e., the number of elements contained in this set). If we had a population of  $n$  unrelated individuals genotyped for  $m$  biallelic loci, then the total number of possible values of  $W$  would be  $3^{nm}$ . However, given the kinship between individuals, the number of possible values of  $W$  is smaller than  $3^{nm}$ . Let  $\mathcal{G}$  be the support of  $\pi(W|P^*)$ , then number of possible values that  $W$  can take is  $|\mathcal{G}|$ , namely the cardinality of the set  $\mathcal{G}$ . To find  $|\mathcal{G}|$ , the pedigree of the population is used because along with the genotypes of founders, it defines how many individuals could potentially have one, two or three genotypes for each marker locus. For example, a progeny from parents with genotypes AA and AA has genotype AA with probability one, while a progeny from parents AA and AB could have genotypes AA or AB with probabilities equal to  $1/2$ . Let  $\mathcal{F}$  be the set of founders, then  $|\mathcal{F}|=f$ , thus there are  $3^{fm}$  possible values for the submatrix of  $W$  corresponding to founders under the assumption that they are unrelated. Hereinafter, each one of these possible values is defined as a “base genotypic configuration”. Notice that each one of these  $fm$  genotypic configurations induces a different set of possible genotypes in the rest of the population. Under base genotypic configuration  $k$ ,  $1 \leq k \leq 3^{mf}$ , for each marker locus the remaining  $n-f$  individuals are grouped into three mutually exclusive sets:  $O_{1j}^k := \{i: |S_{ij} \times D_{ij}|^k = 1, 1 \leq j \leq m, 1 \leq k \leq 3^{mf}\}$ ,  $O_{2j}^k := \{i: |S_{ij} \times D_{ij}|^k = 2, 1 \leq j \leq m, 1 \leq k \leq 3^{mf}\}$ ,  $O_{3j}^k := \{i: |S_{ij} \times D_{ij}|^k = 3, 1 \leq j \leq m, 1 \leq k \leq 3^{mf}\}$ , where  $|S_{ij} \times D_{ij}|^k$  is the cardinality of the set of possible genotypes at marker locus  $j$  resulting from the mating of the parents of individual  $i$  under base genotypic configuration  $k$ ,  $|S_{ij} \times D_{ij}|^k$ . Consequently,  $|O_{lj}^k|$  is the number of individuals in the population for which there are  $l$  possible genotypes at marker  $j$ ,  $1 \leq l \leq 3$  given the  $k^{th}$  base genotypic configuration. Hence, at each marker locus and each base genotypic configuration the following equality is satisfied:  $|O_{1j}^k| + |O_{2j}^k| + |O_{3j}^k| = n - f$ . Therefore, at each marker locus and base genotypic configuration the total number of possible genotypes in the  $n-f$  non-founder individuals is  $|O_{1j}^k| |O_{2j}^k| |O_{3j}^k|$ , and under the linkage equilibrium assumption, the total number of possible genotypes across marker loci given base genotypic configuration  $k$  is

$$\prod_{j=1}^m |O_{1j}^k| |O_{2j}^k| |O_{3j}^k| = \sum_{j=1}^m |O_{2j}^k| \sum_{j=1}^m |O_{3j}^k|$$

Accordingly, given the pedigree of the population, the total number of possible values that matrix  $W$  can take is obtained by summing the above expression over  $k$ :  $|\mathcal{G}| = \sum_{k=1}^{3^{mf}} 2^{\sum_{j=1}^m |O_{2j}^k|} 3^{\sum_{j=1}^m |O_{3j}^k|}$ . As a check of the adequacy of this expression, notice that ignoring pedigree and assuming that all individuals in the population are unrelated is equivalent to treat them all as founders which implies that  $f = n$ , consequently  $|O_{1j}^k| = |O_{2j}^k| = |O_{3j}^k| = 0$ ,  $\forall j = 1, 2, \dots, m$ ,  $\forall k = 1, 2, \dots, 3^{mf}$ , thus  $|\mathcal{G}| = \sum_{k=1}^{3^{mf}} 2^0 3^0 = 3^{nm}$ . Before defining the support of  $W$ , the following sets are defined. The  $k^{th}$  base genotypic configuration is defined as follows:  $\mathcal{G}_{\mathcal{F}}^k := \{w_{ijk} : i \in \mathcal{F}, 1 \leq j \leq m, 1 \leq k \leq 3^{mf}\}$ . For each set  $\mathcal{G}_{\mathcal{F}}^k$ , that is, for each genotypic configuration,  $1 \leq k \leq 3^{mf}$ , define:  $\mathcal{G}_{O_1}^k := \{w_{ij} : i \in O_{1j}^k, 1 \leq j \leq m\}$ ,  $\mathcal{G}_{O_2}^k := \{w_{ij} : i \in O_{2j}^k, 1 \leq j \leq m\}$ ,  $\mathcal{G}_{O_3}^k := \{w_{ij} : i \in O_{3j}^k, 1 \leq j \leq m\}$ . As mentioned before, each set  $\mathcal{G}_{\mathcal{F}}^k$  induces a set  $\mathcal{G}_{O_1}^k \cup \mathcal{G}_{O_2}^k \cup \mathcal{G}_{O_3}^k$ , thus:  $\mathcal{G} = \bigcup_{k=1}^{3^{mf}} \{\mathcal{G}_{\mathcal{F}}^k \cup \mathcal{G}_{O_1}^k \cup \mathcal{G}_{O_2}^k \cup \mathcal{G}_{O_3}^k\}$ .

**Remark 3.** When some individuals are not genotyped or partially genotyped, that is, when a fraction of matrix  $W$  is not observed,  $\pi(W|P^*) = f(W^o|W^N, P^*) \pi(W^N|P^*)$  where  $\pi(W^N|P^*) = \sum_{\mathcal{G}^o} \pi(W|P^*)$ ,  $\mathcal{G}^o$  is the set of possible values of  $W^o$ . However, as will become clear in Section 2.2, explicit computation of  $\pi(W^N|P^*)$  is not required. In this case, some of the elements of  $\pi(W|P^*)$  can be conceptually partitioned as follows:  $n_l^{B_j} = n_{l_o}^{B_j} + n_{l_N}^{B_j}$ ,  $n_l^{A_j} = n_{l_o}^{A_j} + n_{l_N}^{A_j}$ ,  $n^H = n_o^H + n_N^H$  where subindex  $l_o$  indicates that the corresponding count comes from genotyped individuals in the  $l^{th}$  subpopulation and subindex  $l_N$  indicates that the corresponding count comes from non-genotyped individuals.

## 2.2. Full conditionals, homoscedastic residuals, homogeneous and heterogeneous marker effect covariance matrix models

Henceforth, it is assumed that vector  $\mathbf{g}$  and columns of matrix  $W$  are ordered by marker unless otherwise indicated. The full conditionals are denoted as  $\pi(\bullet|Else)$ . Firstly,  $\mathbf{g}|Else \sim MVN \left( \left( I_m \otimes (G^0)^{-1} + \frac{W^o W^o}{\sigma^2} \right)^{-1} \frac{1}{\sigma^2} W^o \mathbf{y}, \left( I_m \otimes (G^0)^{-1} + \frac{W^o W^o}{\sigma^2} \right)^{-1} \right)$ . If  $W_k$  denotes the submatrix of  $W$  corresponding to marker  $k$ ,  $W_k$  is of dimension  $n \times S$  and has the form  $W_k = (w_{1k} \dots w_{nk})'$ ,  $w_{ik} = (0 \dots w_{ik} \dots 0)_{1 \times S}$ ,  $i = 1, 2, \dots, n$ , the only non-null entry of vector  $w_{ik}$  is the random variable corresponding to the genotype of the  $i^{th}$  individual for the  $k^{th}$  marker  $w_{ik}$  and it is located at position  $l$ ,  $l = 1, 2, \dots, S$ , where  $l$  is the subpopulation to which individual  $i$  pertains. Other full conditionals are  $G^0|Else \sim IW \left( a + m, \Sigma + \sum_{j=1}^m \mathbf{g}_j \mathbf{g}_j' \right)$ ,  $\sigma^2|Else \sim IG \left( \frac{v+n}{2}, \frac{(y-W\mathbf{g})(y-W\mathbf{g}) + \tau^2}{2} \right)$ . To arrive at  $\pi(W^N|Else)$  the following definitions have to be made. The rows of  $W$  for individuals with missing genotypes are partitioned as  $W^{Mc}$ ,  $W^{M_1}, \dots, W^{M_K}$  which respectively represent the rows of  $W$  for non-genotyped individuals, and individuals partially genotyped having missing genotypes for loci subsets  $M_1N, \dots, M_KN$ . Accordingly, the subvector of the data vector corresponding to records from non-genotyped or partially genotyped individuals can be partitioned as  $\mathbf{y}^N = (\mathbf{y}^{Mc}, \mathbf{y}^{M_1}, \dots, \mathbf{y}^{M_K})'$ . The rows of  $W$  corresponding to partially genotyped individuals are partitioned as follows:  $W^{M_k o} = (W^{M_k o} : W^{M_k N})$ , where superindex  $M_k o$  denotes the set of loci with observed genotypes, while superindex  $M_k N$  denotes the set of marker loci with missing genotypes. Similarly, when doing computations among these submatrices and  $\mathbf{g}$ , this vector can be arranged as  $(\mathbf{g}^{M_k o} : \mathbf{g}^{M_k N})'$ , then:

$$\pi(W^N|Else) = \pi(W^N | \mathbf{y}^N, W^o, \mathbf{g}, \sigma^2, P^*)$$

$$\propto \pi^+(W|P^*) \exp \left( -\frac{1}{2\sigma^2} (-2\mathbf{g}' W^N \mathbf{y}^N + \mathbf{g}' W^N W^N \mathbf{g}) \right)$$

$$\times \prod_{k=1}^K \exp\left(\frac{-1}{2\sigma^2} h(W^{M_k}, \mathbf{g}^{M_k}, \mathbf{y}^{M_k})\right)$$

where

$$h(W^{M_k}, \mathbf{g}^{M_k}, \mathbf{y}^{M_k}) = 2(\mathbf{g}^{M_k N'} W^{M_k N'} W^{M_k O} \mathbf{g}^{M_k O} - \mathbf{g}^{M_k N'} W^{M_k N'} \mathbf{y}^{M_k}) + \mathbf{g}^{M_k N'} W^{M_k N'} W^{M_k N} \mathbf{g}^{M_k N},$$

$\pi^+(W|P^*) = f^+(W^O|W^N, P^*)\pi(W^N|P^*)$  and  $f^+(W^O|W^N, P^*)$  is the part of the  $W$  component of the likelihood depending on  $W^N$ . Notice that this is a non-standard pmf and that when  $W^O$  depends only on  $W^N$  the form of  $\pi(W^N|Else)$  remains the same because  $f^+(W^O|W^N)\pi(W^N|P^*) = \pi^+(W|P^*)$ . When  $\mathbf{r}$  is known

$$\pi(P|Else) = \pi(P|W^O, \mathbf{r}) = \pi(P|W, \mathbf{r})$$

$$\propto \prod_{j=1}^m P_{(S+1)j}^{\alpha_{S+1}-1} \prod_{l=1}^S \left\{ p_{lj}^{n_l^{B_j} + \alpha_l - 1} (r_l - p_{lj})^{n_l^{A_j}} \right\}$$

which is the product of  $m$  non-standard pdf. Recall that when  $\mathbf{r}$  is unknown, there is a slight difference in this expression as was shown in Section 2.1.

**Remark 4.** In the absence of missing genotypes, that is,  $W^O = W$ , the previous expression is not the full conditional density of  $P$ , but its posterior density.

For the heterogeneous marker effect covariance matrix model  $G$  is a block-diagonal matrix comprised by  $m$  blocks of dimension  $S \times S$  as described in Section 2.1. Under this model  $\pi(G) = \prod_{j=1}^S \pi(G_j)$ . This prior pdf is the only difference with the previous model; therefore, the joint posterior is very similar (see Appendix A). Hence, all full conditionals are the same except for  $\mathbf{g}|Else \sim MVN\left(\left(G^{-1} + \frac{W^O W}{\sigma^2}\right)^{-1} \frac{1}{\sigma^2} W^O \mathbf{y}, \left(G^{-1} + \frac{W^O W}{\sigma^2}\right)^{-1}\right)$ ,  $G^{-1} = \text{Block diag.}(G_j^{-1}), j = 1, 2, \dots, m$  and  $G_j|Else \stackrel{\text{ind}}{\sim} IW(a + 1, \Sigma + \mathbf{g}_j \mathbf{g}_j')$ . The full conditionals for models with heteroscedastic residuals are presented in Appendix A along with joint posteriors.

### 2.3. Model comparison via Deviance Information Criterion

The term null model refers to simplified versions of the proposed models. These null models ignore the factor splitting the complete population into subpopulations; therefore, each marker has a single overall effect and allelic frequencies are assumed to be the same across subpopulations.

Null models are as follows:  $\mathbf{y} = W_0 \mathbf{g}_0 + \boldsymbol{\epsilon}$ , where  $\mathbf{y}$  is the same as before,  $\mathbf{g}_0$  is an  $m \times 1$  unobservable random vector containing allele substitution effects of each marker,  $(W_0)_{n \times m}$  is the random observable design matrix which is of the form  $(W_1' : \dots : W_S')'$  when ordering data by subpopulation, and  $\boldsymbol{\epsilon}$  is a random vector of residuals. The priors for  $\mathbf{g}_0$  are simply univariate versions of the priors used for  $\mathbf{g}$ . Thus,  $\mathbf{g}_0|G^D \sim (\bullet|G^D)$ ,  $G^D = \text{Diag}(\sigma_{g_1}^2, \dots, \sigma_{g_m}^2)$ ,  $\sigma_{g_j}^2 \stackrel{iid}{\sim} IG\left(\frac{a}{2}, \frac{b}{2}\right)$ , (for the homogeneous marker effect variance model  $\sigma_{g_1}^2 = \dots = \sigma_{g_m}^2 = \sigma_g^2$ ) and the residual variance  $\sigma^2$  is given an  $IG\left(\frac{\tau^2}{2}, \frac{v}{2}\right)$  prior as before. In addition,  $\mathbf{p} = (p_1, p_2, \dots, p_m)$  is a vector of overall reference allele frequencies,  $W_0|\mathbf{p} \sim \pi(W_0|\mathbf{p})$  is a simplified version of  $\pi(W|P^*)$  (shown later), and the prior for  $\mathbf{p}$  is  $p_j \stackrel{iid}{\sim} \text{Beta}(\alpha, \beta), j = 1, 2, \dots, m$ .

The Deviance Information Criterion (DIC; Spiegelhalter et al., 2002) combines a measure of goodness of fit based on the posterior distribution and a penalty for model complexity, and despite some criticism it has been used in different areas to perform model comparison (Gelman et al., 2013; Spiegelhalter et al., 2014). It has the following form:

$$DIC = -2 \log f(\text{Data}|\hat{\theta}_B) + 2p_{DIC}$$

where  $p_{DIC} = 2(\log f(\text{Data}|\hat{\theta}_B) - E_{Q[\text{Data}]}[\log f(\text{Data}|\theta)])$ ,  $\hat{\theta}_B = E[\theta|\mathbf{y}]$  is the posterior mean of the unknown parameters. The first component of  $DIC$  is a measure of model adequacy, whereas the second one is the effective number of parameters which is a penalty for increasing model complexity (Spiegelhalter et al., 2002). Models with a smaller DIC are preferred. Recall that for any of our models the likelihood has two components:  $f(\mathbf{y}, W^O|W^N, \mathbf{g}, R, P^*) = f(\mathbf{y}|\mathbf{W}, \mathbf{g}, R)f(W^O|W^N, P^*)$  that were denoted as the  $\mathbf{y}$  component and the  $W$  component. Thus, the general form of the DIC is:

$$DIC = -2 \log f(\mathbf{y}|W^O, \hat{W}_B^N, \hat{\mathbf{g}}_B, \hat{R}_B) + 2p_{DIC-y} - 2 \log f(W^O|\hat{W}_B^N, \hat{P}_B^*) + 2p_{DIC-W}$$

$$= DIC_y + DIC_W$$

where  $p_{DIC-y} = 2(\log f(\mathbf{y}|W^O, \hat{W}_B^N, \hat{\mathbf{g}}_B, \hat{R}_B) - E_{W^N, \mathbf{g}, R, P^*|\mathbf{y}, W^O}[\log f(\mathbf{y}|\mathbf{W}, \mathbf{g}, R)])$  and  $p_{DIC-W} = 2(f(W^O|\hat{W}_B^N, \hat{P}_B^*) - E_{W^N, P^*|\mathbf{y}, W^O}[f(W^O|W^N, P^*)])$ . Thus, as the likelihood, the DIC can be decomposed into a  $\mathbf{y}$  component  $DIC_y$  and a  $W$  component  $DIC_W$ .

### 2.4. Parameter inference via MCMC

In this section, some issues about MCMC algorithms to carry out inference are briefly discussed. Notice that when  $W$  is fully observed, the fact that there are no missing genotypes implies that posterior sampling for the (hyper) parameters of the  $W$  component of the likelihood and the (hyper) parameters of the  $\mathbf{y}$  component can be performed separately. The full conditionals of  $\mathbf{g}, G, \sigma^2, g_0$ , and  $\sigma_g^2$  are known; therefore, samples from the joint posterior can be obtained using a Gibbs sampler (Casella and George, 1992) while samples from the posterior distribution of allelic frequencies can be obtained using a Metropolis-Hastings algorithm. Specifically, independent Metropolis algorithms are considered here. For the scenario of  $\mathbf{r}$  known, the new samples can be generated in two steps: firstly a Dirichlet vector is sampled, and secondly its elements are scaled with the appropriate elements of  $\mathbf{r}$ . Alternatively, uniform(0,1) distributions can be used as proposal, which simplifies computations. With such proposal, given the current state of the chain denoted as  $P^l$ , the acceptance probability of the new sample  $P_+^l$  is  $\min\left\{\frac{\pi(P_+^l|W)}{\pi(P^l|W)}, 1\right\}$ . For null models, the posterior distribution of  $\mathbf{p}_0$  is the product of  $m$  Beta( $p_j; n^{B_j} + \alpha, n^{A_j} + \beta$ ) distributions,  $j = 1, 2, \dots, m$ . Hence, direct sampling can be implemented if needed and the functional form of the posterior mean is known. When  $\mathbf{r}$  is unknown, the candidate to sample from the posterior of  $(p_j, q_j), j = 1, 2, \dots, m$ , could be a Dirichlet of dimension six.

On the other hand, when matrix  $W$  is partially observed a Metropolis-within-Gibbs strategy (Robert and Casella, 2010) can be used to sample from the joint posterior. This strategy is useful due to the fact that nor  $\pi(W^N|Else)$  neither  $\pi(P^*|Else)$  are standard distributions and the existence of the parameter  $W^N$  does not allow to carry out separate sampling algorithms as before because this is a parameter of both components of the likelihood. Accordingly, there are two Metropolis steps in the algorithm to sample from the posterior of the full models. The first one is used to obtain samples from  $\pi(W^N|Else)$ . A good proposal is  $\pi(W^N|W^O, P^*)$  because obtaining direct samples from this distribution via the inverse transform method for discrete random variables (Robert and Casella, 2010) is straightforward. The functional form of  $\pi(W^N|W^O, P^*)$  is derived from first principles as explained in 2.3.1. Thus, given the current state of the chain  $W^{N_l}$ , the acceptance probability of a new sample  $W_+^{N_l}$  is:  $\min\left\{\frac{\pi(W_+^{N_l}|Else)\pi(W_+^{N_l}|W^O, P^*)}{\pi(W^{N_l}|Else)\pi(W^{N_l}|W^O, P^*)}, 1\right\}$ . This applies to both situations:  $\mathbf{r}$  known and  $\mathbf{r}$  unknown. The second Metropolis step is used to draw samples from  $\pi(P|Else)$  for  $\mathbf{r}$  known or  $\pi(P, Q|Else)$  for  $\mathbf{r}$  unknown. The proposals mentioned for the non-missing genotypes scenario also work here. For the null models, it turns out that  $\forall j = 1, 2, \dots, m$ ,  $\pi(p_j|Else)$  is a known distribution, it is a

$\text{Beta}(n^{B_j} + \alpha, n^{A_j} + \beta)$  and consequently only one Metropolis step is needed because direct sampling from the full conditional distribution of  $p_0$  is feasible. Notice that this full conditional distribution is the posterior distribution of  $p_0$  when matrix  $W$  is completely observed.

2.5. Simulation study

In order to provide an example of the implementation of some of the proposed models and the computation of some criteria to compare their performance, two simulated datasets were used. Simulation of these datasets involved two main steps: Simulation of genotypes (QTL and SNP), and simulation of QTL effects and noise. The phenotypes were simulated as the sum of additive genetic effects (sum of QTL allele content times the allele effect) and noise. Datasets were simulated using the software QMSim (Sargolzaei and Schenkel, 2013). In both cases, a historical population was simulated by creating 1000 generations of random mating using a forward-in-time approach in order to reach mutation-drift equilibrium and to create linkage disequilibrium (Sargolzaei and Schenkel, 2013). The historical population size in each generation was 1000 with 500 males and 500 females. Then, subpopulations were created from individuals pertaining to the historical population under different selection pressures and criteria, and different mating systems (Table 1).

Phenotypes were simulated with different number of QTL controlling the trait and different heritabilities. Furthermore, the population structure also differed because the criteria to simulate the subpopulations were different for each trait. Briefly, dataset 1 involved three subpopulations with different number of generations were migration was allowed and the heritability of the trait was high. Dataset 2 comprised two subpopulations with only two generations, no migration and the heritability of the trait was low (Table 1). For further details concerning the simulation see appendix B.

Given that this paper is focused on proposing and explaining a set of across population genome-wide prediction models and not with their large scale implementation, the number of simulated SNP and sample size were low in order to avoid computational issues (Table 1). Phenotype 1 illustrates the situation in which the number of markers is equal to the number of QTL affecting the trait, while for phenotype 2 the number of markers is larger than the number of QTL controlling the trait. These contrasting simulation schemes, different selection pressures and criteria, mating designs and number of generations were

used to mimic real life situations where different subpopulations have different backgrounds. These simulated datasets were used to carry out analyses using the following models: Homogeneous and heterogeneous marker effect covariance matrices with homoscedastic residuals and their null versions. Only models with homoscedastic residuals were used to analyze these datasets because simulations did not consider heteroscedastic residuals.

The analyses performed involved implementation of MCMC algorithms explained in Section 2.4, the computation of DIC and the computation of the following quantities measuring predictive performance and accuracy: the squared correlation between predicted breeding values and phenotypes in the testing populations, hereinafter called predictive ability, and squared correlations between true and predicted breeding values computed in the testing populations (accuracy). Because true breeding values were available for the complete populations, squared correlations between true and predicted breeding values in the training populations were also computed.

For dataset 1, the training population was comprised of generations 0–2 of subpopulation 1, 0–5 from subpopulation 2 and generation 0 of subpopulation 3, while the testing population included generation 3 of subpopulation one, generation 6 of subpopulation 2 and generation 1 of subpopulation 3. For dataset 2, the training population was composed of generations 0 and 1 of subpopulations 1 and 2 and the testing dataset contained generation 2 of subpopulations 1 and 2.

In dataset 2, the full genotypes of three individuals (one founder from each subpopulation and a non-founder from subpopulation 1) were not included in the analysis in order to simulate the case of missing genotypes.

It was assumed that  $r = \left(\frac{1}{S}, \dots, \frac{1}{S}\right)$ . In an initial analysis, a scaled Dirichlet distribution was used as proposal to draw samples from  $\pi(PI|Else)$ , but the behavior of the chains was not satisfactory because the acceptance rate was too low (results not shown). Consequently the product of  $S$  independent uniform  $\left(0, \frac{1}{S}\right)$  distributions was used as proposal. For each dataset, 20.000 iterations were run; the first 10.000 were considered burn-ins. An in-house R script (R Core Team, 2015) was created to carry out the analyses which were performed using the University of Florida’s high performance computing cluster.

3. Results

3.1. Simulated populations

Tables 2 and 3 show features corresponding to characteristics of the simulated genomes and populations.

In both datasets, none of the markers had a minor allele frequency lower than 0.05. Thus, all the simulated marker loci were considered in the analyses.

3.2. DIC, predictive ability and accuracies of predicted breeding values

For dataset 1, the DIC computed using the “W-component” of the likelihood for the full models was 4717671 and 6589105 for the null models. Thus, it provided evidence in favor of the full models when estimating allelic frequencies in the base population. Table 3 shows DIC values for dataset 1, Table 4 DIC values for dataset 2 and Table 5 shows predictive abilities and accuracies in both datasets. For Tables 3–5, the following is the meaning of abbreviations for the different models fitted to datasets 1 and 2:  $M_{IG}$ = full model with Multivariate Gaussian prior and homogeneous marker effect covariance matrices,  $M_{IG}^*$ = full model with Multivariate Gaussian prior and heterogeneous marker effect covariance matrices. Recall that all models assumed homoscedastic residuals. The remaining models with subindex 1 replaced by 0 correspond to null versions of the corresponding full

Table 1  
Parameters and selection criteria to simulate phenotypes.

Parameter	Phenotype 1	Phenotype 2
Heritabilities	0.70, 0.62, 0.54	0.20, 0.15
Phenotypic variances	100, 79, 65	100, 94
Number of QTL	600	40
Number of SNP	600	200
Number of Chromosomes	10	2
Base population structure <sup>a</sup>	1: 28 M, 180F, Phen/L 2: 20 M, 90F, Phen/H 3: 50 M, 500F, Rnd	1: 5 M, 25F, Rnd 2: 20 M, 50F, Phen/H
Number of generations, mating system and selection criteria <sup>b</sup>	1: 3, 0.8, 0.4, As1/Phen, Phen/L 2: 6, 0.7, 0.1, As2/Phen, Phen/H 3: 3, 0.7, 0.2, Rnd, Rnd	1: 2, 1, 0.9, Rnd, Rnd 2: 2, 0.9, 0.3, Rnd, Phen/H

<sup>a</sup> For each line, the first number indicates the subpopulation, items separated by a comma respectively show: number of males, number of females, criterion used to select them (Phen = phenotype, Rnd = random, L = lowest values, H = highest values).  
<sup>b</sup> For each line, the first number indicates the subpopulation, items separated by a comma respectively show: Number of generations, proportion of selected females per generation, proportion of selected males per generation, mating design (Rnd= random, As1=assortative by similarity, As2= assortative by dissimilarity, Phen = phenotype), and selection criterion (same abbreviations as in numeral 2).

**Table 2**

Summary of some characteristics of the simulated populations.

Feature	Dataset 1	Dataset 2
Population size (males, females, total)	883, 1565, 2448	67, 103, 170
Average inbreeding per subpopulation	S1:0.0182, S2: 0.0310, S3:0.0	S1: 0.0 , S2:0.0
Average homozygosity per subpopulation	S1: 0.6240, S2: 0.6359, S3:0.6190	S1:0.6392, S2:0.6283
Phenotype sample mean and SD (in brackets) per subpopulation	S1: -19.78 (13.21) S2: 25.71 (9.60) S3: 0.26 (9.91)	S1: -0.5959 (9.3616) S2:8.9253 (11.9571)

**Table 3**

y component and total DIC for dataset 1.

Model	y component of DIC	Total DIC
$M_{1G}$	33,702.55	4,751,373.55
$M_{1G}^*$	11,599.05	4,729,270.05
$M_{0G}$	15,396.32	6,604,501.32
$M_{0G}^*$	13,008.42	6,602,113.42

**Table 4**

y component, W component and total DIC for dataset 2.

Model	y component of DIC	W component of DIC	Total DIC
$M_{1G}$	1314.0	38,367.4	39,681.4
$M_{1G}^*$	1328.8	38,356.4	39,684.2
$M_{0G}$	1365.6	38,180.3	39,545.9
$M_{0G}^*$	1370.1	38,179.0	39,549.1

**Table 5**

Predictive abilities and accuracies in datasets 1 and 2.

Model	Predictive Ability		Accuracy in testing population		Accuracy in Training population	
	Dataset1	Dataset 2	Dataset1	Dataset2	Dataset1	Dataset2
$M_{1G}$	0.29	0.019	0.27	0.04	0.32	0.17
$M_{1G}^*$	0.76	0.016	0.83	0.03	0.94	0.21
$M_{0G}$	0.53	0.004	0.50	0.07	0.55	0.24
$M_{0G}^*$	0.83	0.013	0.88	0.05	0.88	0.23

models.

Thus, in dataset 1, according to the y component of DIC, for the models with homogeneous marker effect covariance matrices (variances) the null model performed better, while for models with heterogeneous covariance matrices (variances) according to this criterion the full model should be preferred over its null version. When considering the whole likelihood to compute the DIC, the two full models had smaller DIC. Additionally, the model with the smallest DIC, and therefore the “best” one under this criterion was model  $M_{1G}^*$ .

In this dataset the two components of the DIC values and therefore DIC values were similar for all models. The y components of DIC were smaller for the full models. Conversely, the W components were smaller for null models as well as total DIC values.

In dataset 1, according to predictive abilities, the model with the best performance was model  $M_{0G}^*$  while model  $M_{1G}$  had the worst performance. The squared Pearson correlations between true and predicted breeding values in testing dataset 1 suggested that the performance of these models followed a trend similar to that indicated by predictive abilities. In training dataset 1, model  $M_{1G}^*$  yielded the highest accuracy and model  $M_{1G}$  had the smallest accuracy.

Predictive abilities and accuracies in the testing sets were extremely low for dataset 2. Accuracies in training set were higher than those obtained in the testing set; however, they were still low. There were not substantial differences between these squared correlations. Predictive abilities were higher for the full models, while accuracies in testing and training sets were higher for the null models.

## 4. Discussion

### 4.1. General features of the models

A group of hierarchical Bayesian linear regression models to carry out simultaneous genome-wide prediction in several subpopulations accounting for randomness of genotypes was presented. The proposed models differed in the prior distribution assigned to the marker effects and on the assumptions made about residual variances (homogeneous or heterogeneous across subpopulations). The priors for the marker effects were multivariate (univariate) Gaussian and allowed homogeneous or heterogeneous covariance matrices (or variances).

The differences between these models and other regression models currently used in across population genome-wide prediction are: 1) subpopulation-specific effects for each marker are considered and their covariance matrices are modeled explicitly, and 2) genotypes are treated as random variables with a distribution that depends on allelic frequencies as well as on pedigree information. The second feature makes these models different from all other genome-wide prediction models. The distribution of genotypes combines pedigree and genomic information that are not used when randomness of W is ignored. It allows accounting for heterogeneity and correlations of allelic frequencies of the same marker across subpopulations and including individuals with phenotypes and missing genotypes in various loci without carrying out a previous imputation. This is possible because the non-observed part of W, denoted as  $W^N$ , is treated as a parameter and therefore imputation is automatically performed. Another advantage is that the use of a Bayesian approach automatically takes into account uncertainty about the imputed genotypes.

Although most of the paper has been devoted to the models allowing subpopulation-specific effects for each marker (the full models), their univariate versions (the null models) are also contributions of this study. These also allow including individuals with missing genotypes in some or all marker loci without need of external imputation and take into account randomness in genotypes. Therefore, these models could also be used either in single population analyses or to conduct across population genome-wide prediction pooling the data as has been done in previous studies (de Roos et al., 2009; Lund et al., 2011; van den Berg et al., 2015; Wientjes et al., 2015) and was also done here.

Doing a joint analysis has the advantage that the number of phenotypes increases, but in our full models the number of location parameters is also incremented because each marker is allowed to have subpopulation-specific effects; moreover, the number of covariance parameters also increases. The gain in accuracy is achieved when factors such as different QTL effects across subpopulations, differences in linkage phase between QTL and markers, and differences in allelic frequencies and LD patterns make marker effects change substantially from one subpopulation to another. Consequently, the performance of these models may have considerable variation from one dataset to another.

The diagonal blocks of G were assumed to be non-structured. A way reduce dimensionality of the parameter space is to assume certain structure of G. For example, it can be assumed that all covariances and variances are the same, thus, only two parameters per block have to be estimated.

The conditional independence property used to derive  $\pi(WP^*)$  implies that allelic frequencies are estimated in the set of oldest individuals with phenotypes. Here, this set of individuals was referred



to as the base population and individuals pertaining to it were referred to as founders. This was done for pragmatic purposes. However, truncating the pedigree by ignoring individuals without phenotypic records created a group of individuals that may not be the actual base population which is defined as that comprised by ancestors with unknown parents (Henderson, 1974; Kennedy et al., 1988). Conversely, in other cases phenotypic records from this population may be available; thus, estimates of allelic frequencies in the true base population can be obtained. Here, it was further assumed that founders were unrelated which is likely to be false in many situations. However, this assumption has been made in conventional models used to do genetic analysis (Henderson, 1974; Kennedy et al., 1988) because pedigrees are not always completely known. Consequently, what is called the base population is not always the true one. Nevertheless, this assumption seems to be reasonable after so many years of successful artificial selection in animals and plants based on predicted breeding values obtained from these models (Hill, 2014; Gianola and Rosa, 2015).

As discussed in Section 2.1.1, the pmf  $\pi(W|P^*)$  could be derived ignoring pedigree information. Then, this pmf could be found as the product of all  $\pi(w_{ij}^l|p_j^*)$  or the product of binomial distributions for gene content (i.e., the number of copies of the reference allele at each locus) across loci and individuals with each binomial distribution depending on the corresponding allelic frequencies. Notice that this requires reparametrizing the mapping of genotypes, that is, instead of having  $\{-1,0,1\}$  as possible values of an entry of  $W$ , values would be  $\{0,1,2\}$ . In this case, all individuals in the population would be used to estimate allelic frequencies instead of using information from a base population. If pedigree information is available, it can be easily incorporated into the derivation of  $\pi(W|P^*)$  as was shown here and the resulting pmf is not very difficult to evaluate. Furthermore, as mentioned before, direct sampling from this pmf can be done via the inverse transform method for discrete random variables. Notwithstanding, in scenarios where pedigree information is very scarce or not reliable, adding the assumption of independence among individual genotypes and using binomial distributions for the gene content of each individual at each marker locus is an option to model the distribution of matrix  $W$  which would induce a joint pmf similar to those presented in Gianola et al., (2010) and Martínez et al. (2015).

If some individuals with phenotypes have only one known parent, the pmf of their genotypes conditioned on this parent and allelic frequencies can be defined in a similar way as was done in Table 1 for the case of a fully known pedigree (see Appendix C). In this situation, Remark 1 does not hold and the functional form of  $\pi(W|P^*)$  changes which implies that  $\pi(W|Else)$  changes as well.

Regarding assumptions about the distribution of allelic frequencies, our models allow for correlations between them. To do that, priors based on a Dirichlet distribution were used. Using these priors require allelic frequencies to be expressed on a complete population basis. This setting brings parameter  $r$  into the picture. The algebra associated with this parameter is clear and straightforward, but its interpretation may be fuzzy. From an algebraic standpoint, these parameters are upper boundaries posed over allelic frequencies to force them to be in the support of the prior distribution, thus they can be seen as analytic instruments. Nevertheless, their meaning from the population genetics standpoint is not very clear. Perhaps, the easier interpretation when assuming  $r_{1l} = \dots = r_{ml} = r_l$ , is that  $r_l$  is the relative frequency or weight of the  $l^{th}$  subpopulation. However, making claims about the biological interpretation of this set of parameters is beyond the scope of this study.

From a statistical viewpoint, two approaches were proposed. The first one assumed that  $r$  was known (truly known or set to some ad hoc value) and  $r_{1l} = \dots = r_{ml} = r_l$ . In the examples used here all subpopulations were given the same weight, that is,  $r_l = 1/S$ ,  $\forall l = 1, 2, \dots, S$ , a pragmatic decision that has been used in other studies, e.g., Gianola et al. (2010). In this scenario, for all  $j$ ,  $p_j$  is modeled as a scaled Dirichlet vector

which allows non-null covariances between its elements. The second approach assumed that  $r$  was unknown and  $\{r_{ij}\}$  varied across marker loci. For each locus the prior was a Dirichlet over allelic frequencies of both alleles in all subpopulations and it permitted obtaining posterior samples of allelic frequencies and  $r$ . Under the assumption of independence of allelic frequencies, independent priors could be assigned to each marker (e.g., Uniform( $0, r_l$ )) and the validity of this assumption could be tested using criteria as Bayes factors or DIC. If data are pooled and structure is ignored (as done in the null models) the full conditional pdf  $\pi(p_0|Else)$  is known and therefore direct sampling can be implemented when matrix  $W$  is not completely observed. On the other hand, when it is completely observed the posterior of  $p_0$  is known and there is no need of sampling to obtain point estimators. The reason for the full conditional of  $p_0$  being a known distribution but not its posterior in the presence of missing genotypes is that  $W^N$  is an extra parameter in the model and obtaining the marginal posterior of  $p_0$  implies marginalization of  $\pi(W^N, p_0|W^o)$  over  $W^N$  which induces a non-standard pmf.

The derivation of the pmf  $\pi(W|P^*)$  and  $\pi(W_0|p_0)$  not only allow inferences concerning the marker allelic frequencies in the base population, but also allow predictions for non-genotyped or partially genotyped animals without performing a previous imputation. This is likely to increase accuracy of genome-wide predictions because it allows incorporating more phenotypic records. Imputed missing genotypes can be obtained using posterior means or medians of  $W^N$ . However, these outputs have to be viewed as a byproduct because these models were not intended to perform imputation. The imputation of missing genotypes is an underlying process in the prediction of genotypic values of individuals with missing genotypes. Notwithstanding, because samples from the posterior of  $W^N$  are available and computation of imputed genotypes is simple, there could be interest in using this output of the model and in such case the accuracy of the imputation would also be of interest. Hence, although imputation was not a main objective of our models, it is worth making a brief comment on it. Though an assessment of imputation accuracy is a matter for further research, two statements can be made about the imputation process in our models. Firstly, one advantage of the models developed here is that they automatically take into account the uncertainty of imputation (as a consequence of using a Bayesian approach). Conversely, in the standard approach where genotype imputation is the first step and then a random linear regression model is fitted using these imputed values as if they were observations, uncertainty is not taken into account. Secondly, a disadvantage of our models is that they do not incorporate LD information when imputing missing genotypes, a source of information that is used by some of the current imputation methods (Li et al., 2009). Here, pedigree information, phenotypes and allelic frequencies are used for imputation. Thus, benchmarking of the procedure developed here with current and well-accepted procedures is material for future studies. Furthermore, another question that can be addressed in future research is if improving this imputation as discussed later in Section 4.3 has a significant impact on the predictive performance of the models.

As mentioned before, the regression models used in genome-wide prediction treat genotypes as fixed and their effects as random while in the classical quantitative genetics theory genotypes are treated as random and allelic substitution effects as fixed. The set of models developed here are something in between because genotypes are treated as random variables as in classical quantitative genetics, and marker effects are considered random as well like in the standard regression models used in genome-wide prediction. de los Campos et al. (2015b) presented an excellent discussion on the connections between the heritability and the so-called genomic heritability obtained with linear regression models. They show why caution has to be exercised when interpreting the parameters obtained using genomic information due to the fact that sometimes the connection between parameters as the additive genetic variance and the genomic variance

are not straightforward. Similarly, [Gianola et al. \(2015\)](#) discussed the fact that connections between genomic correlations and additive genetic correlations are ambiguous. So far, the Bayesian models proposed in this paper are intended to predict breeding values, phenotypes, and to estimate allelic frequencies in a base population using genomic information and no claim is made about the properties of covariance parameters obtained from them.

The discussion above is relevant because the regression variables are not based on genes, but proxies for the causal variants affecting the phenotypes of interest. However, taking into account these limitations and the high degree of caution needed when interpreting parameters obtained from models using molecular markers, some parameters such as the fraction of additive genetic variance explained by the markers are of interest and our models could be used to estimate these quantities.

The family of models developed here could be applied or adapted to different situations. In the simulation, the case of individuals coming from a common founder population pertaining to subpopulations with different selection criteria and mating systems was considered. Other situations in which this set of models could be useful are: 1) simultaneous evaluation of individuals from different breeds or lines, 2) individuals from the same breed or line performing under different environmental conditions (e.g., different geographic regions, production systems, etc.), 3) a combination of numerals 1 and 2, 4) simultaneous evaluation of several correlated traits. In this last case, if all individuals have records for all phenotypes, the design matrix satisfies  $W = I_S \otimes W_+$ , where  $W_+$  is the matrix of dimension  $n \times m$  containing genotypes of  $n$  individuals at  $m$  marker loci. In this case the model is being adapted to handle correlations between the effects of a given marker locus for different traits in a single population. Consequently, for a given choice of prior and assumption about residuals (heteroscedastic or homoscedastic) the model involves the corresponding hierarchical structure except for the pmf of  $W$  conditional on the allelic frequencies and pedigree which is  $\pi(W_+ | p_0^*)$  instead of  $\pi(W | p^*)$ . Recent studies have developed Bayesian multiple-trait genome-wide regression models and have shown that predictions from them are more accurate than those coming from genomic univariate models ([Jia and Jannink, 2012](#)). The hierarchical Bayesian multivariate genome-wide prediction models proposed by [Jia and Jannink \(2012\)](#) have similar components to the models presented here such as the priors for  $g$ , but they do not account for randomness of genotypes. Another step to accommodate our models for multiple-trait prediction is to allow correlated residuals, that is, a non-diagonal matrix  $R$ . In this case, an inverse Wishart prior can be assigned instead of the inverse gamma prior used here.

#### 4.2. Simulation results

As stated in [Section 2.4](#), the aim of this limited simulation was to provide an illustration of the implementation of models and methods developed in this study. Thus, results are not conclusive and further research involving analyses based on more elaborate simulations as well as real datasets to have a better evaluation of the performance of this family of models is needed. Nevertheless, some insights and comments derived from the analyses of these two datasets can be discussed.

The correlation between phenotypes and predicted breeding values (or its square) is one of the most widely used measurements to compare genome-wide prediction models, it is associated with the response to selection and it is easy to compute. On the other hand, as mentioned previously, the DIC combines measures of model adequacy and complexity ([Spiegelhalter et al., 2002](#)).

For dataset 1, the squared correlation between phenotypes and predicted breeding values (the predictive ability) did not show an advantage in predictive capability of models taking into account the population structure, i.e., the existence of the subpopulations ([Table 5](#)). While measures based on squared correlations did not provide

conclusive evidence in favor of the full models, the DIC favored the full models.

As expected, the predictive ability and the other correlations were much smaller in dataset 2 due to the lower heritability of the trait. Although all predictive abilities were low, according to this criterion the performance of the full models was slightly better. Accuracies of predicted breeding values suggested a tiny superiority of null models. The two subpopulations simulated in this dataset diverged by just two generations which could cause only small differences in allelic frequencies, this scenario clearly favors the null models. Accordingly, the DIC component coming from genotypes was slightly better (smaller) for null models as opposed to the case of dataset 1. The total DIC gave evidence in favor of null models. Among predictive ability, accuracy and DIC, accuracy and DIC favored the null models, but the values were very close. The performance of the fitted models was more similar in this dataset than in dataset 1.

In our small simulations, when subpopulations diverged by several generations, migration was allowed and heritabilities were high (dataset 1), full models had better performance in terms of DIC. Conversely, when populations diverged by only a few generations, there was no migration, and heritabilities were low (dataset 2) null models tended to perform better according to this criterion. However, the differences were small. On the other hand, predictive abilities showed a different pattern. In dataset 1 this criterion was higher for null models while in dataset 2 it was smaller for null models. Another feature shown by these simulations was the high variability in model performance that may exist among populations. In dataset 1, according to all criteria except the  $W$  component of DIC, the performance of model  $M_{IG}$  tended to be remarkably poorer while this was not the case in dataset 2.

Other authors have found modest or null increments in predictive performance of models allowing heterogeneous marker effects across subpopulations compared to pooling data and analyzing the complete population as a single one ([Olson et al., 2012](#); [Makgahlela et al., 2013](#); [de los Campos et al., 2015a](#)). All the aforementioned studies used real data from plants and animals. Working with three plant populations and using a model very similar to those proposed here, [Lehermeier et al. \(2015\)](#) found cases in which the strategy of pooling data and ignoring structure performed better and other cases where multivariate models yielded better predictive performance. These authors found that in highly differentiated populations within group and multivariate analyses performed better while the converse occurred in closely related subpopulations with small sample sizes. Roughly speaking, these results are in agreement with the results found in this study.

Using predictive ability, [Lund et al. \(2011\)](#) found a higher accuracy of predicted additive breeding values when pooling the data compared with individual analyses. Similar results were found by [de Roos et al. \(2009\)](#) when heritability was low, divergence of populations was small (small number of generations) and marker density was high (more persistent phase), and by [Wientjes et al. \(2015\)](#) when the QTL effects did not change across subpopulations. Pooling data and ignoring the population structure corresponds to the null models defined in this study, except that models considered by the authors just cited did not account for randomness of genotypes. In our simulation, individual analyses were not considered. Sample size is one of the factors affecting the accuracy of genome-wide predictions ([Meuwissen et al., 2001](#); [Goddard, 2009](#); [Zhong et al., 2009](#)). Presumably it was one of the leading factors causing the results found by [Lund et al. \(2011\)](#). In addition, the Holstein breed is highly inbred and there were several individuals connecting the different populations; this probably made them similar. On the other hand, the studies of [de Roos et al. \(2009\)](#) and [Wientjes et al. \(2015\)](#) used simulated data and explored different scenarios. Both studies found situations in which pooling data was not advantageous.

### 4.3. Refinements and extensions

In this section, some comments regarding possible extensions and refinements of different aspects of the family of models presented in the study are briefly discussed.

In the derivation of the joint pmf of  $W$  conditional on  $P^*$  and pedigree information, row-wise dependence due to kinship was taken into account by using pedigree information to accommodate relationships among genotypes of related individuals. This task was highly simplified due to the conditional independence argument that permitted to find a simpler decomposition of the joint pmf and therefore, a simpler algebraic expression. However, the possible existence of column-wise dependence due to LD was ignored here in order to make the problem more tractable from the mathematical point of view. This is an assumption frequently used in theoretical studies in quantitative genetics and it is well-accepted at least in studies concerned with first approximations to a given problem. For example, Gianola et al. (2009) treated a series of theoretical aspects of some of the Bayesian regression models used in genome-wide prediction using the assumption of linkage equilibrium which implies the mutual independence of the columns of  $W$  used here (they also developed some results accounting for LD in the Appendix). Most of the models currently used in genome-wide prediction are also based on this assumption, few approximations to deal with consequences of LD have been proposed (Gianola et al., 2003; Yang and Tempelman, 2012), but these have not yet been adopted in routine genetic evaluations. Their models do not consider randomness in the genotypes; thus, a consequence of considering LD in these models is the need to account for covariances between marker effects at different loci. Consequently, a refinement of our family of models in this regard, would be to accommodate LD, which can be performed at two levels: 1) account for correlations among columns of  $W$ , and 2) use a non-block-diagonal  $G$  matrix.

A potential consequence of accounting for non-independence of the columns of  $W$  could be the reduction in the cardinality of  $\mathcal{G}$  that is induced by the fact that the number of possible values of a column of  $W$  depends on the values at one or more different columns (as it happened with rows). Another assumption made here was the absence of mutations which caused that when conditioning on the genotypes of the parents of an individual, the probabilities of its genotype taking a given value were completely defined by the parental genotypes, making this random variable conditionally independent of allelic frequencies. Thus, another refinement in  $\pi(W|P^*)$  would be to account for mutation. Therefore, the derivation of  $\pi(W|P^*)$  to accommodate dependence between columns of  $W$  and mutation, and the impact of this refinement on predictive performance and the accuracy of imputed genotypes (if it is of interest) pose a problem for further research.

If relationships among founders (as defined in this paper) were to be taken into account, from the theoretical point of view it is not hard to visualize how to do it. For the sake of simplicity, the case of two individuals and one locus is considered; consequently, the sub-index associated with locus is omitted. Let  $W_1$ ,  $W_2$  be the genotypes of individuals 1 and 2, and  $W_C$  the genotypes of the set of relevant common ancestors. Suppose that 1 is not a parent of 2. Then:

$$\begin{aligned}\pi(W_1, W_2|P^*) &= \sum_{\mathcal{G}^C} \pi(W_1, W_2|W_C, P^*) \pi(W_C|P^*) \\ &= \sum_{\mathcal{G}^C} \pi(W_1|W_C, P^*) \pi(W_2|W_C, P^*) \pi(W_C|P^*),\end{aligned}$$

where  $\mathcal{G}^C$  is the set of possible values that the set of genotypes of relevant common ancestors can take according to the pedigree (as explained in Section 2.1.1) and the second equality follows from the conditional independence of the genotypes of individuals 1 and 2 given the common ancestors and allelic frequencies. By relevant common ancestors it is meant that the genotypes of these ancestors provide information about the genotypes of 1 and 2 when conditioning on the

full set of common ancestors, i.e., if  $\mathcal{D}$  is the whole set of common ancestors then  $\mathcal{D} = \mathcal{C} \cup \mathcal{C}^c$  (the super-index  $c$  means complement with respect to  $\mathcal{D}$ ) and  $\pi(W_1, W_2|W_{\mathcal{D}}, P^*) = \pi(W_1, W_2|W_C, P^*)$ . Notice that unless individuals 1 and 2 are full sibs, their conditional pmf given the relevant common ancestors depends on  $P^*$ . Of course, it makes  $\pi(W|P^*)$  a more complex expression and reduces the cardinality of  $\mathcal{G}$ . See Appendix D for a toy example of  $\pi(W_1, W_2|W_C, P^*)$  when 1 and 2 are half sibs. Although the problem is tractable from the theoretical standpoint, it may be difficult to compute these values especially with complex pedigrees where the set of common ancestors may be large such as those found in animal and plant populations. The example in Appendix D shows that even in a simple case, computation of  $\pi(W_1, W_2|P^*)$  is involved.

## 5. Conclusions

The main contribution of this paper is the theoretical development of a set of models for across population genome-wide prediction incorporating marker genotypes not only as explanatory variables of regression models, but also as realizations of random variables providing information about allelic frequencies and missing genotypes. Although models were intended for across population analysis, they can also be applied in single population studies and adapted for multiple-trait prediction.

Theoretical and computational issues along with possible applications as well as some extensions and refinements of these models pose several problems for future research. Our models treat both genotypes and marker allelic substitution effects as random; therefore, they combine features from classical quantitative genetics theory and traditional genome-wide prediction models.

Some features of the models developed in this study make them promising for genome-wide prediction. Among these, the ability to include phenotypes from individuals with missing genotypes at some or all loci without the need of previous imputation and accounting for uncertainty about imputed genotypes as well as heterogeneity of allelic frequencies across subpopulations are perhaps the most appealing. Further research to assess their performance and also to compare them with other models used in genome-wide prediction is needed.

## Acknowledgments

Authors acknowledge Dr. Malay Ghosh from the Department of Statistics of the University of Florida for useful comments and discussions, and for pointing out relevant references. C. A. Martínez also thanks PhD students Hunter Merrill and Isaac Duerr, and Dr. Nikolay Bliznyuk from the Department of Agricultural and Biological Engineering of the University of Florida for their advice in computational issues, Fulbright Colombia and “Departamento Administrativo de Ciencia, Tecnología e Innovación” COLCIENCIAS for supporting his PhD and Master programs at the University of Florida through a scholarship, and Bibiana Coy for her love, support and constant encouragement.

## Appendices. Supporting information

Supplementary data associated with this article can be found in the online version at <http://dx.doi.org/10.1016/j.jtbi.2016.12.020>.

## References

- Bernardo, R., Yu, J., 2007. Prospects for genomewide selection for quantitative traits in maize. *Crop Sci.* 47, 1082–1090.
- Casella, G., George, E.I., 1992. Explaining the Gibbs Sampler. *Am. Stat.* 46 (3), 167–174.
- Casella, G., Berger, R., 2002. *Statistical Inference* 2nd ed. Duxbury, Pacific Grove, CA, USA.
- Chen, L., Li, C., Miller, S., Schenkel, F., 2014. Multi-population genomic prediction using a multi-task Bayesian learning model. *BMC Genet.* 15, 53.

- Core Team, R., 2015. R: a language and environment for statistical computing. R foundation for statistical computing, Vienna, Austria (<https://www.R-project.org/>).
- de los Campos, G., Gianola, D., Allison, D.B., 2010. Predicting genetic predisposition in humans: the promise of whole-genome markers. *Nat. Rev. Genet.* 11, 880–886.
- de los Campos, G., Sorensen, D., Gianola, D., 2015b. Genomic heritability: what is it? *PLoS Genet.* 11 (5), e1005048.
- de los Campos, G., Veturri, Y., Vázquez, A.I., Lehermeier, C., Pérez-Rodríguez, P., 2015a. Incorporating genetic heterogeneity in whole-genome regressions using interactions. *J. Agric., Biol. Environ. Stat.* 20 (4), 467–490.
- Desta, Z.A., Ortiz, R., 2014. Genomic selection: genome-wide prediction in plant improvement. *Trends Plant Sci.* 19 (9), 592–601.
- Falconer, D.S., Mackay, T.F.C., 1996. *Introduction to Quantitative Genetics* 4th ed.. Longmans Green, Harlow, UK.
- Gelman, A., Carlin, J.B., Stern, H., Dunson, D.B., Vehtari, A., Rubin, D.B., 2013. *Bayesian Data Analysis* 3rd ed. Chapman and Hall/CRC, Boca Raton, FL, USA.
- Gianola, D., 2013. Priors in whole-genome regression: the Bayesian alphabet returns. *Genetics* 194, 573–596.
- Gianola, D., Rosa, G., 2015. One hundred years of statistical developments in animal breeding. *Annu. Rev. Anim. Biosci.* 3, 19–56.
- 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., Simianer, H., Qanbari, S., 2010. A two-step method for detecting selection signatures using genetic markers. *Genet. Res. Camb.* 92 (2), 141–155.
- Gianola, D., de los Campos, G., Hill, W.G., Manfredi, E., Fernando, R.L., 2009. Additive genetic variability and the Bayesian alphabet. *Genetics* 183, 347–363.
- Gianola, D., de los Campos, G., Toro, M.A., Naya, H., Schön, C.C., Sorensen, D., 2015. Do molecular markers inform about pleiotropy? *Genetics* 201, 23–29.
- Goddard, M.E., 2009. Genomic selection: prediction of accuracy and maximization of long term response. *Genetica* 136, 245–257.
- Goddard, M.E., Hayes, B.J., 2007. Genomic Selection. *J. Anim. Breed. Genet.* 124, 323–330.
- Guttmacher, A.E., Collins, F.S., 2002. Genomic medicine- a primer. *New Engl. J. Med.* 347, 1512–1520.
- Hayes, B.J., Bowman, P.J., Chamberlain, A.J., Goddard, M.E., 2009. Invited review: genomic selection in dairy cattle: progress and challenges. *J. Dairy Sci.* 92, 433–443.
- Henderson, C.R., 1974. Use of all relatives in intraherd prediction of breeding values and producing abilities. *J. Daity Sci.* 58 (12), 1910–1916.
- Hill, W.G., 1984. On selection among groups with heterogeneous variance. *Anim. Prod.* 39 (3), 473–477.
- Hill, W.H., 2014. Applications of population genetics to animal breeding, from wright, fisher and lush to genomic prediction. *Genetics* 196, 1–16.
- Huang, H., Windig, J.J., Vereijken, A., Calus, M.P.L., 2014. Genomic prediction based on data from three layer lines using non-linear regression models. *Genet. Sel. Evol.* 46, 75.
- Jia, Y., Jannink, J.L., 2012. Multiple-trait genomic selection methods increase genetic value prediction accuracy. *Genetics* 192, 1513–1522.
- Karoui, S., Carabaño, M.J., Díaz, C., Legarra, A., 2012. Joint genomic evaluation of French dairy cattle breeds using multiple-trait models. *Genet. Sel. Evol.* 44, 39.
- Kennedy, B.W., Schaeffer, L.R., Sorensen, D.A., 1988. Genetic properties of animal models. *J. Dairy Sci.* 71 (2), 17–26.
- Lehermeier, C., Schon, C., de los Campos, G., 2015. Assessment of genetic heterogeneity in structured plant populations using multivariate whole-genome regression models. *Genetics* 201, 323–337.
- Li, Y., Willer, C., Sanna, S., Abecasis, G., 2009. Genotype imputation. *Annu. Rev. Genom. Hum. Genet.* 10, 387–406.
- Lund, M.S., de Roos, A.P.W., de Vries, A.G., Druet, T., Ducrocq, V., Fritz, S., Guillaume, F., Gulbrandtsen, B., Liu, Z., Reents, R., Schrooten, C., Seefried, F., Su, J., 2011. A common reference population from four European Holstein populations increases reliability of genomic predictions. *Genet. Sel. Evol.* 43, 43.
- Lynch, M., Walsh, E., 1998. *Genetics and Analysis of Quantitative Traits*. Sinauer associates Inc, Sunderland, MA, USA.
- Makgahlela, M.L., Mantysaari, E.A., Strandén, I., Koivula, M., Nielsen, U.S., Sillanpää, M.J., Juga, J., 2013. Across breed multi-trait random regression genomic predictions in the Nordic Red dairy cattle. *J. Anim. Breed. Genet.* 130, 10–19.
- Martínez, C.A., Khare, K., Elzo, M.A., 2015. On the Bayesness, minimaxity and admissibility of point estimators of allelic frequencies. *J. Theor. Biol.* 383, 106–115.
- 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.
- Olson, K.M., VanRaden, P.M., Tooker, M.E., 2012. Multibreed genomic evaluations using purebred Holsteins, Jerseys, and Brown Swiss. *J. Dairy Sci.* 95, 5378–5383.
- Robert, C.P., Casella, G., 2010. *Introducing Monte Carlo Methods with R*. Springer, New York, NY, USA.
- de Roos, A.P.W., Hayes, B.J., Goddard, M.E., 2009. Reliability of genomic predictions across multiple populations. *Genetics* 183, 1545–1553.
- Sargolzaei, M., Schenkel, F.S., 2013. *QMSim User's Guide Version 1.10*. Centre for Genetic Improvement of Livestock, Department of Animal and Poultry Science. University of Guelph, Guelph, Canada.
- Spiegelhalter, D.J., Best, N.G., Carlin, B.P., van der Linde, A., 2002. Bayesian measures of model complexity and fit (with discussion). *J. R. Stat. Soc. Ser. B* 64, 583–639.
- Spiegelhalter, D.J., Best, N.G., Carlin, B.P., van der Linde, A., 2014. The deviance information criterion: 12 years on. *J. R. Stat. Soc. Ser. B* 76, 485–493.
- van den Berg, S., Calus, M.P.L., Meuwissen, T.H.E., Wientjes, Y.C.J., 2015. Across population genomic prediction scenarios in which Bayesian variable selection outperforms GBLUP. *BMC Genet.* 16, 416.
- Wientjes, Y.C.J., Veerkamp, R.F., Bijma, P., Bovenhuis, H., Schrooten, C., Calus, M.P.L., 2015. Empirical and deterministic accuracies of across-population genomic prediction. *Genet. Sel. Evol.* 47, 5.
- Wright, S., 1930. Evolution in Mendelian populations. *Genetics* 16, 98–159.
- Wright, S., 1937. The distribution of genetic frequencies in populations. *Genetics* 23, 307–320.
- Yang, W., Tempelman, R.J., 2012. A Bayesian antedependence model for whole genome prediction. *Genetics* 190, 1491–1501.
- Zhong, S., Dekkers, J.C.M., Fernando, R.L., Jannink, J.K., 2009. Factors affecting accuracy from genomic selection in populations derived from multiple inbred lines: a Barley case study. *Genetics* 182, 355–364.