1 **Joint genome-wide prediction in several populations accounting for randomness of genotypes:**

2 **A hierarchical Bayes approach. I: Multivariate Gaussian priors for marker effects and**

3 **derivation of the joint probability mass function of genotypes**

4 Carlos Alberto Martínez[a,b], Kshitij Khare[b], Arunava Banerjee[c], Mauricio A. Elzo[a]

5 [a]Department of Animal Sciences

6 [b]Department of Statistics

7 [c]Department of Computer and Information Science and Engineering

8 University of Florida, Gainesville, FL, USA

9 Correspondence: Carlos Alberto Martínez, Department of Animal Sciences, University of Florida,

10 Gainesville, FL 32611, USA.

11 Tel: 352-328-1624.

12 Fax: 352-392-7851.

13 E-mail: carlosmn@ufl.edu

14

15

16

17

18

19

20

21

22

23

24

25

**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.

Key words: Across population genome-enabled prediction; Bayesian modeling; heterogeneous allelic frequencies; distribution of genotypes.

**1. Introduction**

50　The use of molecular markers located across the whole genome for prediction of breeding values

51　(Meuwissen et al., 2001) and phenotypes (Goddard and Hayes, 2007, Gianola et al., 2009) has

52　proven to be a useful tool in animals (Hayes et al., 2009), humans (Guttmacher et al., 2002; de los

53　Campos et al., 2010) and plants (Bernardo and Yu, 2007; Desta and Ortiz, 2014). This success has

54　given rise to a tremendous amount of research in the area of statistical genomics in order to obtain

55　better genome-wide predictions (Goddard and Hayes, 2007; Gianola, 2013; Hill, 2014; Gianola and

56　Rosa, 2015).

57　Most of the methods have been developed for prediction in a single population. Across population

58　studies usually use predictions obtained from individual populations or pool data to perform a single

59　analysis (de Roos et al., 2009). On one hand, pooling data and performing a single analysis may

60　increase the accuracy of genome-wide prediction because the number of records has an important

61　impact on it (Meuwissen et al., 2001; Goddard, 2009; Zhong et al., 2009). On the other hand, it may

62　decrease accuracy when the effects of QTL controlling the trait are not the same across populations

63　(de Roos et al., 2009; van den Berg et al., 2015; Wientjes et al., 2015).

64　Analyzing data from Holstein cattle performing in different European countries, Lund et al. (2011)

65　reported that pooling data and carrying out a single analysis increased the accuracy of genomic

66　predictions. With simulated data, de Roos et al. (2009) found that pooling data was beneficial when

67　populations had diverged by few generations, marker density was high and heritability was low, but

68　for more distant populations and less dense marker panels they found a small decrease in accuracy.

69　Using simulated data, Wientjes et al. (2015) studied the effect of differences in QTL allele

70　substitution effects across populations on the accuracy of genome-wide prediction. They found that

71　when allele substitution effects changed across populations, the accuracies decreased in proportion to

72　the genetic correlation between populations. Using the same dataset, van den Berg et al. (2015)

73　looked for across population genomic prediction scenarios under which Bayesian variable selection

74　models had a better performance than genomic BLUP (GBLUP). They concluded that Bayesian

75 variable selection models outperform GBLUP when the number of QTL is small as in single

76 population analyses, but the difference in accuracy is larger in the across population case.

77 None of these studies allowed marker effects to differ from one population to another. However, de

78 Roos et al. (2009) highlighted the need for alternative methods that allow population-specific

79 estimation of allele substitution effects in across population genome wide prediction. Chen et al.

80 (2014) proposed a Bayesian model with different SNP effects for each population that permits

81 sharing information across populations through a common set of latent variables indicating weather a

82 given marker is associated with a QTL or not. They did not model covariance matrices of marker

83 effects explicitly. With real and simulated data they found that this model increased the accuracy of

84 across population genome-wide prediction, especially when the number of QTL was small and

85 correlations among QTL effects from different populations were high. Recently, Bayesian models

86 that account for genetic heterogeneity have been proposed. Multivariate models considering

87 correlated population specific marker effects were developed by Lehermeir et al. (2015) while de los

88 Campos et al. (2015a) proposed a model with main marker effects and interactions. Using real data

89 from three plant populations, Lehermeir et al. (2015) found cases in which the strategy of pooling

90 data and ignoring structure performed better and others where the multivariate models yielded better

91 predictive performance. For example, in highly differentiated populations within group and

92 multivariate analyses performed better. Using real datasets from pigs and wheat, de los Campos et al.

93 (2015a) found modest superiority of the interaction model relative to the model using pooled data

94 and the model that analyzed each subpopulation separately. Similar studies have implemented

95 multivariate models in multibreed dairy cattle populations (Karoui et al., 2012; Olson et al., 2012;

96 Makgahlela et al., 2013). Huang et al. (2014) used non-linear models to perform genome wide

97 prediction in layer hens when the reference population was comprised by individuals from several

98 breeds or lines and compared them with a multiple-trait GBLUP model. They found that the various

99 models used had a similar predictive performance.

4

100  If several populations are to be evaluated simultaneously, the possible existence of genotype by

101  environment interaction, lack of persistence of linkage phase and variation in allelic frequencies

102  across populations indicate the need for an analysis that accounts for the fact that combining them

103  creates a structured complete population. It has been reported that population structure may act as an

104  effect modifier (de los Campos et al., 2015a). Furthermore, it has to be considered that not only the

105  allele substitution effects of a particular locus in different populations may be correlated, but also its

106  frequencies in each population (e.g., due to gene flow).

107  Another feature that has been overlooked in the random linear regression models used in genome-

108  wide prediction is the randomness of the matrix containing a one to one mapping from the set of

109  genotypes to a subset of the integers, namely the design matrix. This matrix is treated as fixed in

110  genome-wide prediction models, while in classical quantitative genetics theory it is treated as random

111  (Falconer and Mackay, 1996; Lynch and Walsh, 1998). Besides being in agreement with the classical

112  theory, taking into account the randomness of this matrix, that is, the randomness of genotypes,

113  permits the estimation of allelic frequencies because when treated as an observable discrete random

114  matrix, its probability mass function (pmf) depends on the allelic frequencies. Thus, under a

115  Bayesian setting, allelic frequencies are treated as random because these are unknown parameters.

116  Further, the works of Wright (1930; 1937) provide additional support to treat allelic frequencies as

117  random variables making Bayesian inference even more attractive.

118  Thus, the objective of this study was to propose hierarchical Bayesian models to carry out

119  simultaneous genome-wide prediction in several populations accounting for randomness of marker

120  genotypes, heterogeneity and correlation of allelic frequencies across populations, and population-

121  specific allelic substitution effects.

122  **2. Methods**

123  *2.1 The models*

124    Hereinafter the complete population or simply the population is defined as the set of individuals with

125    phenotypes considered in the study. Suppose that there exists some criterion (e.g., environment, race,

126    breed, line, etc.) to split this population into $\mathcal{S}$ subpopulations. To make the problem more tractable,

127    some simplifying assumptions are made. The first one is linkage equilibrium. The second one is

128    Hardy-Weinberg equilibrium. The third one is that starting from the oldest individuals with

129    phenotypes, the pedigree is fully known. Lastly, mutations are ignored.

130    The basic linear model used to describe the relationship between response variables and marker allele

131    substitution effects is $\boldsymbol{y} = W\boldsymbol{g} + \boldsymbol{e}$, where $\boldsymbol{y}$ is a vector containing dependent variables (e.g., records

132    corrected for non-genetic factors), $W$ is an observable random matrix containing a one to one

133    mapping from individual marker genotypes to a subset of the integers to be defined later, $\boldsymbol{g}$ is an

134    unknown random vector of marker allelic substitution effects for every population and $\boldsymbol{e}$ is a random

135    vector of residuals. A more detailed notation is the following. If records are sorted by subpopulation

136    as well as the columns of $W$ and the elements of $\boldsymbol{g}$, then for every $l = 1,2,\dots,\mathcal{S}$, $\boldsymbol{y}_l = W_l\boldsymbol{g}_l + \boldsymbol{e}_l$,

137    with dimensions: $(\boldsymbol{y}_l)_{n_l \times 1}$, $(W_l)_{n_l \times m}$, $(\boldsymbol{g}_l)_{m \times 1}$ and $(\boldsymbol{e}_l)_{n_l \times 1}$ where $n_l$ is the sample size of

138    subpopulation $l$, and $m$ is the number of marker loci. Thus, the total sample size is $n = \sum_{l=1}^{\mathcal{S}} n_l$.

139    The scenario where only a part of matrix $W$ is observed because some individuals are not genotyped

140    or individuals are genotyped for different numbers of marker loci is also considered. This is done by

141    treating this non-observed part of $W$ as a parameter in the model as it will be explained later.

142    The case of diploid individuals and biallelic marker loci is considered. The effect of every marker

143    locus is defined as the regression of records on a function of the number of copies of the reference

144    allele and in quantitative genetics it corresponds to the allele substitution effect (Falconer and

145    Mackay, 1996; Lynch and Walsh, 1998). The number of copies can be "centered" at zero giving the

146    following codification. Let A and B be the marker alleles at each locus and let B be the reference

147    allele. Then:

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

148    Different versions of the hierarchy that represents the stochastic component of each model were

149    considered. Models vary according to the assumptions on the variance of residuals and the priors

150    posed over the marker effects. The most parsimonious model is the one considering homoscedastic

151    residuals and homogeneous marker effect covariance matrices. The hierarchical Bayesian model

152    assuming homoscedastic residuals and multivariate Gaussian priors for marker effects has the

153    following structure:

$$\boldsymbol{y}|W, \boldsymbol{g}, \sigma^2 \sim MVN(W\boldsymbol{g}, \sigma^2 I)$$

$$W|\boldsymbol{p}_1^*, \boldsymbol{p}_2^*, \dots, \boldsymbol{p}_m^* \sim \pi(\cdot|\boldsymbol{p}_1^*, \boldsymbol{p}_2^*, \dots, \boldsymbol{p}_m^*)$$

$$\boldsymbol{p}_j^* \overset{iid}{\sim} \pi(\boldsymbol{p}^*), j = 1,2, \dots, m$$

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

$$\boldsymbol{g}|G \sim MVN(0, G), G = Block\ Diag\ \{G_j\}_{j=1}^{m}$$

$$G_j \overset{iid}{\sim} Inverse\ Wishart(a, \boldsymbol{\Sigma}) := IW(a, \boldsymbol{\Sigma})$$

$$G_j = \begin{bmatrix} \sigma_{j_1}^2 & \sigma_{j_{1,2}} & \cdots & \sigma_{j_{1,S}} \\ & \sigma_{j_2}^2 & \cdots & \sigma_{j_{2,S}} \\ & & \ddots & \vdots \\ sym & & & \sigma_{j_S}^2 \end{bmatrix}$$

154    where $\sigma^2$ is the residual variance, $\sigma_{j_l}^2$ is the variance of the effect of the $j^{th}$ marker in the $l^{th}$

155    subpopulation, $\sigma_{j_{l,l'}}$ is the covariance between effects of marker $j$ in subpopulations $l$ and $l'$, $\boldsymbol{p}_j^*$ is a

156    parameter associated with allelic frequencies of the $j^{th}$ marker in each subpopulation and $\pi(\boldsymbol{p}^*)$ is its

157    density. Details on these parameters and their probability density function (pdf) are given later.

158  In the case of heterogeneous residual variances across subpopulations, residual variances $\sigma_1^2, \ldots, \sigma_S^2$

159  are given independent $IG\left(\frac{\tau^2}{2}, \frac{v}{2}\right)$ priors and then: $\boldsymbol{y}|W, \boldsymbol{g}, R \sim MVN(W\boldsymbol{g}, V)$, $R = (\sigma_{e1}^2, \ldots, \sigma_{eS}^2)$ and

160  $V = Block\ Diag.\left\{\sigma_{el}^2 I_{n_l}\right\}_{l=1}^{S}$. Hill (1984) found that in the presence of heterogeneous

161  environmental variances, across population analyses assuming homogenous residuals variances

162  yielded an excess of individuals selected from populations with higher environmental variances. This

163  is why heterogeneity of residual variances across subpopulations was considered in this study.

164  The general framework assumes that in each subpopulation there is a fraction of genotyped

165  individuals and a fraction of non-genotyped or partially genotyped individuals. Let $W^\sigma$ and $W^N$

166  denote the observed (data) and non-observed (an unknown parameter) parts of $W$. Let $P^* =$

167  $(\boldsymbol{p}_1^*, \boldsymbol{p}_2^*, \ldots, \boldsymbol{p}_m^*)$;  therefore,  $\pi(W|P^*) = \pi(W^\sigma, W^N|P^*)$  can  be  expressed  as:

168  $f(W^\sigma|W^N, P^*)\pi(W^N|P^*)$. Thus, the full likelihood has the form:

$$f(\boldsymbol{y}, W^\sigma|W^N, \boldsymbol{g}, R, P^*) = f(\boldsymbol{y}|W^\sigma, W^N, \boldsymbol{g}, R, P^*)f(W^\sigma|W^N, \boldsymbol{g}, R, P^*)$$

$$= f(\boldsymbol{y}|W, \boldsymbol{g}, R)f(W^\sigma|W^N, P^*).$$

169  Henceforth, $f(\boldsymbol{y}|W, \boldsymbol{g}, R)$ will be referred to as the $\boldsymbol{y}$ component of the likelihood and

170  $f(W^\sigma|W^N, P^*)$ will be referred to as the $W$ component.

171  The simplest case for the covariance matrix of marker effects is $G = I \otimes G^0$. Under this setting the

172  assumption is that the covariance structure is the same for all markers. This is statistically convenient

173  due to the fact that the number of covariance parameters is reduced. Further, in analysis considering a

174  single population, it has been found that specifying a different variance for each marker does not

175  allow too much Bayesian learning about marker effect variances (Gianola et al., 2009). Here, models

176  assigning the same covariance matrix to the effects of all marker loci and models considering a

177  different covariance matrix for the effects of each marker locus were considered and these models

178  were referred to as homogeneous marker effect covariance matrix models and heterogeneous marker

179    effect covariance matrix models. Let $\mathcal{P}_S^+$ denote the space of symmetric positive definite matrices of

180    dimension $S \times S$. Then, the marginal prior distribution of $\boldsymbol{g}$ is:

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

181    For details, see Appendix A. Similarly, for the heterogeneous marker effect covariance matrix model

182    it can be shown (appendix A) that: $\pi(\boldsymbol{g}) \propto \dfrac{1}{\prod_{j=1}^m \left(1 + \frac{1}{a+1-S}\boldsymbol{g}_j'\boldsymbol{\Sigma}_*^{-1}\boldsymbol{g}_j\right)^{\left(\frac{a+1}{2}\right)}}$, which is the product of $m$

183    multivariate t distributions with scale matrix $\boldsymbol{\Sigma}_* = \frac{1}{a+1-S}\boldsymbol{\Sigma}$ and degrees of freedom $a + 1 - S$;

184    therefore, under this prior, marker effects are marginally independent and identically distributed. At

185    this point, the following remark can be made.

186    *Remark 1* Under the assumption of homogeneous marker effect covariance matrices, *a priori* the

187    marker effects are marginally dependent. This happens because when integrating with respect to the

188    common covariance matrix $G^0$, the term $\sum_{j=1}^m \boldsymbol{g}_j\boldsymbol{g}_j'$ and the hyper-hyperparameter $\boldsymbol{\Sigma}$ are factored,

189    resulting in a function that cannot be written as the product of $m$ functions, each one depending on a

190    different $\boldsymbol{g}_j$. Moreover, the joint prior density is not standard.

191    To take into account the belief that allelic frequencies of the same marker vary across subpopulations

192    and may be correlated, the prior $\pi(\boldsymbol{p}^*)$ is built based on a Dirichlet distribution. To do that, the allelic

193    frequency of the reference allele in marker locus $j$ in subpopulation $l$ has to be expressed on a

194    complete population basis, that is, $p_{lj}$ is expressing the frequency of the reference allele in locus $j$ in

195    subpopulation $l$ relative not to subpopulation $l$, but to the complete population. Thus, the frequencies

196    of the two alleles at a given marker locus and a given subpopulation do not add to one, but to some

197    sort of relative frequency of that subpopulation in that locus denoted as $r_{lj}$. Let $\boldsymbol{r} = (\boldsymbol{r}_1, \dots, \boldsymbol{r}_S), \boldsymbol{r}_l =$

198    $(r_{l1}, \dots, r_{lm}), l = 1,2, \dots, S$. With this parameterization $\sum_{l=1}^S p_{lj} \leq 1, \forall\, j = 1,2, \dots, m$, with equality if

199    and only if the reference allele is fixed in all subpopulations. Conversely, allelic frequencies

200    expressed on a subpopulation basis satisfy the constraint that the sum of the frequencies of the two

201    alleles at each marker locus equals one within each subpopulation. Let $q_{jl}, j = 1,2, \ldots, m, l =$

202    $1,2, \ldots, \mathcal{S}$, be the frequencies of the non-reference alleles expressed on a complete population basis,

203    then $p_{lj} + q_{lj} = r_{lj}$. The two parameterizations of allelic frequencies are related by the one to one

204    mapping $p_{lj}^* = p_{lj}/r_{lj}$.

205    Consider the case when $\boldsymbol{r}$ is known and $r_{l1} = \cdots = r_{lm} = r_l \ \forall \ l$. Then, elements of vector $\boldsymbol{r} =$

206    $(r_1, \ldots, r_{\mathcal{S}})$ can be seen as subpopulation weights, that is, they are related to subpopulation sizes. By $\boldsymbol{r}$

207    being known, it is meant that it is either actually known or it is specified following some assumption.

208    A pragmatic decision would be to assign equal subpopulation weights, an assumption that was

209    also made in other studies (e.g., Gianola et al. 2010). Once $\boldsymbol{r}$ has been specified, there is an extra

210    restriction over each $\boldsymbol{p}_j = (p_{1j}, \ldots, p_{\mathcal{S}j})$. For $l = 1,2, \ldots, \mathcal{S}$ the following condition must be satisfied:

211    $p_{lj} \leq r_l$. Therefore, the support of the distribution of $\boldsymbol{p}_j$ given $\boldsymbol{r}$ is $\Omega_j^r := \{\boldsymbol{p}_j \in \mathbb{R}^{\mathcal{S}} | 0 < p_{lj} \leq$

212    $r_l \ \forall \ l, \sum_{l=1}^{\mathcal{S}} r_l = 1\}$. Notice that the condition $\sum_{l=1}^{\mathcal{S}} r_l = 1$ implies that vectors in $\Omega_j^r$ satisfy $\sum_{l=1}^{\mathcal{S}} p_{lj} \leq$

213    1. Thus, under this approach the prior used for each $\boldsymbol{p}_j$ is one corresponding to a scaled Dirichlet

214    random vector. If $\boldsymbol{\beta} = (\beta_1, \ldots, \beta_{\mathcal{S}}) \sim Dirichlet(\boldsymbol{\alpha})$, $\boldsymbol{\alpha} \in \mathbb{R}^{\mathcal{S}+1}$, then the prior assigned to $\boldsymbol{p}_j$ is the

215    distribution of vector $(\beta_1 r_1, \ldots, \beta_{\mathcal{S}} r_{\mathcal{S}})$ which clearly pertains to $\Omega_j^r$. Then, the pdf $\pi(\boldsymbol{p}_j | \boldsymbol{r})$ is derived

216    using standard results from the theory of distributions of transformations of random variables

217    (Casella and Berger, 2002). This derivation is simplified by the fact that the transformation is linear

218    and therefore the Jacobian is constant. It follows that: $\pi(\boldsymbol{p}_j | \boldsymbol{r}) \propto \prod_{l=1}^{\mathcal{S}} \left\{\left(\frac{p_{lj}}{r_l}\right)^{\alpha_l - 1}\right\} p_{(\mathcal{S}+1)j}^{\alpha_{\mathcal{S}+1} - 1}$, where

219    $p_{(\mathcal{S}+1)j} = 1 - \sum_{l=1}^{\mathcal{S}} \frac{p_{lj}}{r_l}$.

220    The second approach is to assume that $\boldsymbol{r}$ is unknown. The density $\pi(\boldsymbol{p} | \boldsymbol{r})$ could be used and a

221    Dirichlet distribution could be assigned to each $\boldsymbol{r}_j$ adding one more level to the hierarchy. However,

222    using $p_{lj} + q_{lj} = r_{lj}$ and properties of the Dirichlet distribution, the following strategy allows

223    assigning a prior to allelic frequencies and the weights $r$ without putting an extra level in the

224    hierarchy. To this end it is assumed that $r_{lj}$ varies for each $j$ and each $l$. A $Dirichlet\left(\left(\boldsymbol{\alpha}_p, \boldsymbol{\alpha}_q\right)\right)$

225    prior is posed over $\left(\boldsymbol{p}_j, \boldsymbol{q}_j\right)$, where $\boldsymbol{q}_j$ is the analog of $\boldsymbol{p}_j$ for the non-reference allele at each locus

226    and $\boldsymbol{\alpha}_p = \left(\alpha_{1p}, \dots, \alpha_{Sp}\right), \boldsymbol{\alpha}_q = \left(\alpha_{1q}, \dots, \alpha_{Sq}\right)$. Consequently, by properties of the Dirichlet

227    distribution it follows that $\boldsymbol{r}_j \sim Dirichlet\left(\left(\alpha_{1p} + \alpha_{1q}, \dots, \alpha_{Sp} + \alpha_{Sq}\right)\right)$.

228    *2.1.1 Deriving the joint pmf of marker genotypes conditional on allelic frequencies*

229    Given the kinship structure of a population (i.e., the pedigree) one can find several generations

230    comprised of genotyped, partially genotyped and non-genotyped individuals. Therefore, the approach

231    is to derive the pmf of the complete matrix $W$, i.e., the joint pmf of individuals with phenotypic

232    records. Under this setting, $m$ is the total number of marker loci to be included in the analysis (it

233    usually corresponds to the size of the densest marker panel used in the population).

234    Across columns, that is, across marker loci, the problem is simplified by assuming linkage

235    equilibrium, which implies independence of genotypes at different loci. Therefore, for an arbitrary

236    subpopulation, the joint density of its column vectors is simply the product of their marginal pmf.

237    When considering all subpopulations, the same assumption implies that marker genotypes at different

238    loci are independent. The following derivations hold for any of the previously discussed approaches

239    to model allelic frequencies distributions.   Under the assumption of Hardy-Weinberg equilibrium it

240    follows that marginally:

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

241    Recall that $p_{lj}^* = p_{lj}/r_{lj}$. Notice that $p_{lj}^*$ is used instead of $p_{lj}$ because it allows defining a proper

242    pmf in the sense that the sum of the probabilities of the three possible values of $w_{ij}^l$ equals one

243    (which does not happen when using $p_{lj}$). The pmf $\pi\left(w_{ij}^l \mid p_{lj}^*\right)$ can be also written as:

11

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

244 where $I_{zi}$ is the indicator variable of the mutually exclusive events $w_{ij}^l = z, z \in \{-1,0,1\}$. By the

245 linkage equilibrium assumption it follows that for individual $i$ in population $l$: $\pi(\boldsymbol{w}_i^l|\boldsymbol{p}_j^*) =$

246 $\prod_{j=1}^m \pi(w_{ij}^l|p_{lj}^*)$.

247 The rows of matrix $W$ represent individuals with records. Because of the kinship between them, the

248 genotype of a given individual is not independent of the genotype of their relatives. Furthermore, this

249 non-independence has to be considered across subpopulations (e.g., half or full sibs may pertain to

250 different subpopulations). This approach is based on the pedigree of the complete population. The

251 "base" animals or "founders" can be pragmatically defined as the oldest individuals with phenotypic

252 records and those individuals with phenotypes and unknown parents. To facilitate computations, it is

253 assumed that these individuals are unrelated. Hereinafter this set is referred to as the base population,

254 and individuals in this set are referred to as founders or base individuals. The remaining individuals

255 in the population are referred to as non-founders. This pmf could be derived ignoring pedigree

256 information which is equivalent to mutual independence of the rows of $W$, then $\pi(W|P^*) =$

257 $\prod_{j=1}^m \prod_{l=1}^S \prod_{i=1}^{n_l} \pi(w_{ij}^l|p_{lj}^*)$. However, this would ignore information contained in the pedigree and

258 would unnecessarily make the parametric space of $W^N$ larger, which does not seem to be the best

259 way to proceed.

260 The ordering of individuals is arbitrary, but a convenient way to do it here is according to the

261 pedigree in such a way that the founders are given the first indices. For marker locus $j$ in population $l$

262 the target is to find:

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

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

$$\pi(\boldsymbol{w}_j^l|p_{lj}^*) = \pi(w_{n_l j}^l|w_{1j}^l, \dots, w_{(n_l-1)j}^l, p_{lj}^*)\pi(w_{1j}^l, \dots, w_{(n_l-1)j}^l|p_{lj}^*)$$

$$= \pi\left(w_{n_l j}^l \middle| w_{1j}^l, \ldots, w_{(n_l-1)j}^l, p_{lj}^*\right) \pi\left(w_{(n_l-1)j}^l \middle| w_{1j}^l, \ldots, w_{(n_l-2)j}^l, p_{lj}^*\right) \pi\left(w_{1j}^l, \ldots, w_{(n_l-2)j}^l \middle| p_{lj}^*\right)$$

$$= \pi\left(w_{n_l j}^l \middle| w_{1j}^l, \ldots, w_{(n_l-1)j}^l, p_{lj}^*\right) \cdots \pi\left(w_{1j}^l \middle| p_{lj}^*\right)$$

$$= \prod_{i=0}^{n_l-2} \left\{ \pi\left(w_{(n_l-i)j}^l \middle| w_{1j}^l, \ldots, w_{(n_l-i-1)j}^l, p_{lj}^*\right) \right\} \pi\left(w_{1j}^l \middle| p_{lj}^*\right).$$

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

$$\pi(W^l | \boldsymbol{p}_l^*) = \prod_{i=0}^{n_l-2} \left\{ \pi\left(\boldsymbol{w}_{n_l-i}^l \middle| \boldsymbol{w}_1^l, \ldots, \boldsymbol{w}_{n_l-i-1}^l, \boldsymbol{p}_l^*\right) \right\} \pi\left(\boldsymbol{w}_1^l \middle| \boldsymbol{p}_l^*\right),$$

265   where each one of the pmf $\pi\left(\boldsymbol{w}_{n_l-i}^l, \middle| \boldsymbol{w}_1^l, \ldots, \boldsymbol{w}_{n_l-i-1}^l, \boldsymbol{p}_l^*\right)$ is the product:

266   $\prod_{j=1}^m \pi\left(w_{(n_l-i)j}^l, \middle| w_{1j}^l, \ldots, w_{(n_l-i-1)j}^l, p_{lj}^*\right), 0 \leq i \leq n_l - 2$ and $\pi\left(\boldsymbol{w}_1^l \middle| \boldsymbol{p}_l^*\right) = \prod_{j=1}^m \pi\left(w_{1j}^l \middle| p_{lj}^*\right).$

267   Now, a conditional independence argument is used to simplify $\pi(W^l | \boldsymbol{p}_l^*)$. Given the genotypes of the

268   parents of individual $i$, its genotype is independent of the genotype of collateral relatives and other

269   ancestors. It is possible that the parents of individual $i$ in population $l$ pertain to subpopulations $l^*$

270   and $l'$. Thus, at this point the complete population is considered. In addition, notice that given the

271   parental genotypes, the genotype of an individual does not depend on the allelic frequencies because

272   this conditional pmf is determined using basic segregation rules (see Appendix A). From these

273   arguments it follows that for individual $i$, $\pi(\boldsymbol{w}_i | \boldsymbol{w}_1, \ldots, \boldsymbol{w}_{i-1}, P^*) = \pi\left(\boldsymbol{w}_i | \boldsymbol{w}_{S_i}, \boldsymbol{w}_{D_i}\right)$, where $\boldsymbol{w}_{S_i}$ and

274   $\boldsymbol{w}_{D_i}$ are the genotypes of the parents of individual $i$. The pmf of non-founder genotypes at marker

275   locus $j$ conditioned on their parental genotypes is presented in Appendix A. Therefore, $\pi(W|P^*)$ can

276   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

277   considering the rows corresponding to founders and $W_{NF}$ is the submatrix of $W$ comprised of the

278   rows corresponding to non-founders. Let $f$ be the total number of founders. Under the assumption

279   that these individuals are unrelated, the pmf of their genotypes given allelic frequencies is:

13

$$\pi(W_F|P^*) = \prod_{i=1}^{f} \pi(W_i|P^*) = \prod_{j=1}^{m}\prod_{i=1}^{f} \pi(w_{ij}|P^*) = \prod_{j=1}^{m}\prod_{l=1}^{S}\prod_{i=1}^{f_l} \pi(w_{ij}^l|p_{lj}^*)$$

$$= \prod_{j=1}^{m}\prod_{l=1}^{S}\prod_{i=1}^{f_l} (p_{lj}^{*2})^{I_{1i}} (2p_{lj}^*(1-p_{lj}^*))^{I_{0i}} ((1-p_{lj}^*)^2)^{I_{-1i}}$$

$$= \prod_{j=1}^{m}\prod_{l=1}^{S} (p_{lj}^{*2})^{n_l^{BB_j}} (2p_{lj}^*(1-p_{lj}^*))^{n_l^{AB_j}} ((1-p_{lj}^*)^2)^{n_l^{AA_j}}$$

$$= \prod_{j=1}^{m}\prod_{l=1}^{S} 2^{n_l^{AB_j}} p_{lj}^{*\, 2n_l^{BB_j}+n_l^{AB_j}} (1-p_{lj}^*)^{2n_l^{AA_j}+n_l^{AB_j}} = 2^{n^H} \prod_{j=1}^{m}\prod_{l=1}^{S} p_{lj}^{*\, n_l^{B_j}} (1-p_{lj}^*)^{n_l^{A_j}},$$

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

$$\pi(W_F|P,\boldsymbol{r}) = 2^{n^H} \prod_{j=1}^{m}\prod_{l=1}^{S} \frac{1}{r_{lj}^{2f_l}} p_{lj}^{\, n_l^{B_j}} (r_{lj}-p_{lj})^{n_l^{A_j}}$$

281    where $f_l$ is the number of founders in the $l^{th}$ subpopulation; thus, $f = \sum_{l=1}^{S} f_l$, $n_l^{BB_j}, n_l^{AB_j}$ and $n_l^{AA_j}$

282    are the counts of founders with genotypes BB, AB and AA at marker locus $j$ in subpopulation $l$

283    respectively, $n_l^{B_j} = 2n_l^{BB_j} + n_l^{AB_j}$ is the total count of B alleles at marker locus $j$ in founders from

284    subpopulation $l$, $n_l^{A_j} = 2n_l^{AA_j} + n_l^{AB_j}$ is the total count of A alleles at marker locus $j$ in founders

285    from subpopulation $l$ and $n^H = \sum_{j=1}^{m}\sum_{l=1}^{S} n_l^{AB_j}$ is the total number of heterozygous loci in the base

286    population. In terms of the random variables $w_{ij}^l$, $n_l^{BB_j}, n_l^{AB_j}$ and $n_l^{AA_j}$ can be written as: $n_l^{BB_j} =$

287    $\sum_{i=1}^{f_l} I_{1i}$, $n_l^{AA_j} = \sum_{i=1}^{f_l} I_{-1i}$, $n_l^{AB_j} = f_l - \left(n_l^{BB_j} + n_l^{AA_j}\right) = f_l - \sum_{i=1}^{f_l}(w_{ij}^l)^2$.

288    For non-founders:

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

289    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

290    subpopulation $l$. Hence:

14

$$\pi(W|P^*) = \prod_{j=1}^{m}\prod_{l=1}^{S}\prod_{i=1}^{f_l} \pi(w_{ij}^l|p_{lj}^*) \times \prod_{j=1}^{m}\prod_{l=1}^{S}\prod_{i'=f_l+1}^{n_l} \pi\left(w_{i'j}|w_{S_{i'j}}, w_{D_{i'j}}\right)$$

$$= \prod_{j=1}^{m}\prod_{l=1}^{S}\prod_{i=1}^{f_l} \left\{ \pi(w_{ij}^l|p_{lj}^*) \prod_{i'=f_l+1}^{n_l} \pi\left(w_{i'j}^l|w_{S_{i'j}}, w_{D_{i'j}}\right) \right\}$$

$$= 2^{n^H} \prod_{j=1}^{m}\prod_{l=1}^{S} \left\{ p_{lj}^{* n_l^{B_j}} (1 - p_{lj}^*)^{n_l^{A_j}} \prod_{i'=f_l+1}^{n_l} \pi\left(w_{i'j}^l|w_{S_{i'j}}, w_{D_{i'j}}\right) \right\}$$

$$\Rightarrow \pi(W|P, \boldsymbol{r}) = 2^{n^H} \prod_{j=1}^{m}\prod_{l=1}^{S} \left\{ \frac{1}{r_{lj}^{2f_l}} p_{lj}^{n_l^{B_j}} (r_{lj} - p_{lj})^{n_l^{A_j}} \prod_{i'=f_l+1}^{n_l} \pi\left(w_{i'j}^l|w_{S_{i'j}}, w_{D_{i'j}}\right) \right\}.$$

291    *Remark 2* Under the assumptions presented at the beginning of this section, given base genotypes, the

292    process defining the inheritance of alleles is completely determined by the pedigree information. The

293    pedigree allows tracing the set of possible values that genotypes can take from a given individual

294    back to the base population. It implies that allelic frequencies have to be known only in the base

295    population because the distribution of genotypes in the set of non-founders is completely determined

296    by the pedigree. Stated another way, given the pedigree, only the founder genotypes carry

297    information about allelic frequencies.

298    The next step is to formally define the support (set of values of $W$ with non-null probability) of the

299    pmf $\pi(W|P^*)$ and its cardinality (i.e., the number of elements contained in this set). If we had a

300    population of $n$ unrelated individuals genotyped for $m$ biallelic loci, then the total number of possible

301    values of $W$ would be $3^{nm}$. However, given the kinship between individuals, the number of possible

302    values of $W$ is smaller than $3^{nm}$. Let $\mathcal{G}$ be the support of $\pi(W|P^*)$, then number of possible values

303    that $W$ can take is $|\mathcal{G}|$, namely the cardinality of the set $\mathcal{G}$. To find $|\mathcal{G}|$, the pedigree of the population

304    is used because along with the genotypes of founders, it defines how many individuals could

305    potentially have one, two or three genotypes for each marker locus. For example, a progeny from

306    parents with genotypes AA and AA has genotype AA with probability one, while a progeny from

15

307     parents AA and AB could have genotypes AA or AB with probabilities equal to ½. Let $\mathcal{F}$ be the set

308     of founders, then $|\mathcal{F}| = f$, thus there are $3^{fm}$ possible values for the submatrix of $W$ corresponding

309     to founders under the assumption that they are unrelated. Hereinafter, each one of these possible

310     values is defined as a "base genotypic configuration". Notice that each one of these $fm$ genotypic

311     configurations induces a different set of possible genotypes in the rest of the population. Under base

312     genotypic configuration $k, 1 \leq k \leq 3^{mf}$, for each marker locus the remaining $n - f$ individuals are

313     grouped into three mutually exclusive sets: $O_{1j}^k := \left\{ i \colon \left| \{S_{ij} \times D_{ij}\}^k \right| = 1, 1 \leq j \leq m, 1 \leq k \leq 3^{mf} \right\}$,

314     $O_{2j}^k := \left\{ i \colon \left| \{S_{ij} \times D_{ij}\}^k \right| = 2, 1 \leq j \leq m, 1 \leq k \leq 3^{mf} \right\}$,         $O_{3j}^k := \left\{ i \colon \left| \{S_{ij} \times D_{ij}\}^k \right| = 3, 1 \leq j \leq \right.$

315     $\left. m, 1 \leq k \leq 3^{mf} \right\}$, where $\left| \{S_{ij} \times D_{ij}\}^k \right|$ is the cardinality of the set of possible genotypes at marker

316     locus $j$ resulting from the mating of the parents of individual $i$ under base genotypic configuration

317     $k, \{S_{ij} \times D_{ij}\}^k$. Consequently, $\left| O_{lj}^k \right|$ is the number of individuals in the population for which there are

318     $l$ possible genotypes at marker $j, 1 \leq l \leq 3$ given the $k^{th}$ base genotypic configuration. Hence, at

319     each marker locus and each base genotypic configuration the following equality is satisfied: $\left| O_{1j}^k \right| +$

320     $\left| O_{2j}^k \right| + \left| O_{3j}^k \right| = n - f$. Therefore, at each marker locus and base genotypic configuration the total

321     number of possible genotypes in the $n - f$ non-founder individuals is $1^{\left| O_{1j}^k \right|} 2^{\left| O_{2j}^k \right|} 3^{\left| O_{3j}^k \right|}$, and under the

322     linkage equilibrium assumption, the total number of possible genotypes across marker loci given

323     base genotypic configuration $k$ is

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

324     Accordingly, given the pedigree of the population, the total number of possible values that matrix $W$

325     can take is obtained by summing the above expression over $k$: $|\mathcal{G}| = \sum_{k=1}^{3^{mf}} 2^{\sum_{j=1}^{m} \left| O_{2j}^k \right|} 3^{\sum_{j=1}^{m} \left| O_{3j}^k \right|}$. As a

326     check of the adequacy of this expression, notice that ignoring pedigree and assuming that all

327     individuals in the population are unrelated is equivalent to treat them all as founders which implies

328    that $f = n$, consequently $|O_{1j}^k| = |O_{2j}^k| = |O_{3j}^k| = 0, \forall j = 1,2, \dots, m, \forall k = 1,2, \dots, 3^{mn}$, thus

329    $|\mathcal{G}| = \sum_{k=1}^{3^{mn}} 2^0 3^0 = 3^{mn}$. Before defining the support of $W$, the following sets are defined. The $k^{th}$

330    base genotypic configuration is defined as follows: $\mathcal{G}_{\mathcal{F}}^k := \{w_{ijk} : i \in \mathcal{F}, 1 \leq j \leq m, 1 \leq k \leq 3^{mf}\}$.

331    For each set $\mathcal{G}_{\mathcal{F}}^k$, that is, for each genotypic configuration, $1 \leq k \leq 3^{mf}$, define: $\mathcal{G}_{O_1}^k := \{w_{ij} : i \in$

332    $O_{1j}^k, 1 \leq j \leq m\}$, $\mathcal{G}_{O_2}^k := \{w_{ij} : i \in O_{2j}^k, 1 \leq j \leq m\}$, $\mathcal{G}_{O_3}^k := \{w_{ij} : i \in O_{3j}^k, 1 \leq j \leq m\}$. As mentioned

333    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: $\mathcal{G} = \bigcup_{k=1}^{3^{mf}} \{\mathcal{G}_{\mathcal{F}}^k \cup \mathcal{G}_{O_1}^k \cup \mathcal{G}_{O_2}^k \cup \mathcal{G}_{O_3}^k\}$.

334    *Remark 3* When some individuals are not genotyped or partially genotyped, that is, when a fraction

335    of matrix $W$ is not observed, $\pi(W|P^*) = f(W^o|W^N, P^*)\pi(W^N|P^*)$ where $\pi(W^N|P^*) =$

336    $\sum_{\mathcal{G}^o} \pi(W|P^*)$, $\mathcal{G}^o$ is the set of possible values of $W^o$. However, as will become clear in section 2.2,

337    explicit computation of $\pi(W^N|P^*)$ is not required. In this case, some of the elements of $\pi(W|P^*)$

338    can be conceptually partitioned as follows: $n_l^{Bj} = n_{l_o}^{Bj} + n_{l_N}^{Bj}, n_l^{Aj} = n_{l_o}^{Aj} + n_{l_N}^{Aj}, n^H = n_o^H + n_N^H$

339    where subindex $l_o$ indicates that the corresponding count comes from genotyped individuals in the

340    $l^{th}$ subpopulation and subindex $l_N$ indicates that the corresponding count comes from non-genotyped

341    individuals.

342    *2.2 Full conditionals, homoscedastic residuals, homogeneous and heterogeneous marker effect*

343    *covariance matrix models*

344    Henceforth, it is assumed that vector $\boldsymbol{g}$ and columns of matrix $W$ are ordered by marker unless

345    otherwise indicated. The full conditionals are denoted as $\pi(\cdot|Else)$. Firstly,

346    $\boldsymbol{g}|Else \sim MVN\left(\left(I_m \otimes (G^0)^{-1} + \frac{w'w}{\sigma^2}\right)^{-1} \frac{1}{\sigma^2} W'\boldsymbol{y}, \left(I_m \otimes (G^0)^{-1} + \frac{w'w}{\sigma^2}\right)^{-1}\right)$. If $W_k$ denotes the

347    submatrix of $W$ corresponding to marker $k$, $W_k$ is of dimension $n \times S$ and has the form $W_k =$

348    $(\boldsymbol{w}_{1k}' \quad \cdots \quad \boldsymbol{w}_{nk}')', \boldsymbol{w}_{ik} = (0 \quad \cdots \quad w_{ik} \quad \cdots \quad 0)_{1 \times S}, i = 1,2, \dots, n$, the only non-null entry of vector

349    $\boldsymbol{w}_{ik}$ is the random variable corresponding to the genotype of the $i^{th}$ individual for the $k^{th}$ marker $w_{ik}$

350    and it is located at position $l, l = 1,2, \dots, S$, where $l$ is the subpopulation to which individual $i$

17

351  pertains.      Other      full      conditionals      are      $G^0|Else \sim IW\left(a + m, \Sigma + \sum_{j=1}^m \boldsymbol{g}_j \boldsymbol{g}_j'\right)$,

352  $\sigma^2|Else \sim IG\left(\frac{v+n}{2}, \frac{(\boldsymbol{y}-W\boldsymbol{g})'(\boldsymbol{y}-W\boldsymbol{g})+\tau^2}{2}\right)$. To arrive at $\pi(W^N|Else)$ the following definitions have to be

353  made. The rows of $W$ for individuals with missing genotypes are partitioned as $W^{M_C}, W^{M_1}, \ldots., W^{M_K}$

354  which respectively represent the rows of $W$ for non-genotyped individuals, and individuals partially

355  genotyped having missing genotypes for loci subsets $M_1 N, \ldots, M_K N$. Accordingly, the subvector of

356  the data vector corresponding to records from non-genotyped or partially genotyped individuals can

357  be partitioned as $\boldsymbol{y}^N = (\boldsymbol{y}^{M_C\prime}, \boldsymbol{y}^{M_1\prime}, \ldots., \boldsymbol{y}^{M_K\prime})'$. The rows of $W$ corresponding to partially genotyped

358  individuals are partitioned as follows: $W^{M_k} = (W^{M_k\sigma} : W^{M_k N})$, where superindex $M_k\sigma$ denotes the

359  set of loci with observed genotypes, while superindex $M_k N$ denotes the set of marker loci with

360  missing genotypes. Similarly, when doing computations among these submatrices and $\boldsymbol{g}$, this vector

361  can be arranged as $(\boldsymbol{g}^{M_k\sigma\prime} : \boldsymbol{g}^{M_k N\prime})'$, then:

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

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

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

362  where

$$h(W^{M_k}, \boldsymbol{g}^{M_k}, \boldsymbol{y}^{M_k}) = 2(\boldsymbol{g}^{M_k N\prime}W^{M_k N\prime}W^{M_k\sigma}\boldsymbol{g}^{M_k\sigma} - \boldsymbol{g}^{M_k N\prime}W^{M_k N\prime}\boldsymbol{y}^{M_k}) + \boldsymbol{g}^{M_k N\prime}W^{M_k N\prime}W^{M_k N}\boldsymbol{g}^{M_k N},$$

363  $\pi^+(W|P^*) = f^+(W^\sigma|W^N, P^*)\pi(W^N|P^*)$ and $f^+(W^\sigma|W^N, P^*)$ is the part of the $W$ component of

364  the likelihood depending on $W^N$. Notice that this is a non-standard pmf and that when $W^\sigma$ depends

365  only on $W^N$ the form of $\pi(W^N|Else)$ remains the same because $f^+(W^\sigma|W^N)\pi(W^N|P^*) =$

366  $\pi^+(W|P^*)$. When $\boldsymbol{r}$ is known

$$\pi(P|Else) = \pi(P|W^\sigma, W^N, \boldsymbol{r}) = \pi(P|W, \boldsymbol{r})$$

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

367    which is the product of $m$ non-standard pdf. Recall that when $r$ is unknown, there is a slight

368    difference in this expression as was shown in section 2.1.

369    *Remark 4* In the absence of missing genotypes, that is, $W^\sigma = W$, the previous expression is not the

370    full conditional density of $P$, but its posterior density.

371    For the heterogeneous marker effect covariance matrix model $G$ is a block-diagonal matrix

372    comprised by $m$ blocks of dimension $S \times S$ as described in section 2.1. Under this model $\pi(G) =$

373    $\prod_{l=1}^{S} \pi(G_j)$. This prior pdf is the only difference with the previous model; therefore, the joint

374    posterior is very similar (see Appendix A). Hence, all full conditionals are the same except for

375    $\boldsymbol{g}|Else \sim MVN\left(\left(G^{-1} + \frac{W'W}{\sigma^2}\right)^{-1} \frac{1}{\sigma^2} W'\boldsymbol{y}, \left(G^{-1} + \frac{W'W}{\sigma^2}\right)^{-1}\right), G^{-1} = Block\ diag.\left(G_j^{-1}\right), j = 1,2 \dots, m$

376    and $\quad G_j|Else \overset{ind}{\sim} IW\left(a + 1, \boldsymbol{\Sigma} + \boldsymbol{g}_j\boldsymbol{g}_j'\right)$. The full conditionals for models with heteroscedastic

377    residuals are presented in Appendix A along with joint posteriors.

378    *2.3 Model comparison via Deviance Information Criterion*

379    The term null model refers to simplified versions of the proposed models. These null models ignore

380    the factor splitting the complete population into subpopulations; therefore, each marker has a single

381    overall effect and allelic frequencies are assumed to be the same across subpopulations.

382    Null models are as follows: $\boldsymbol{y} = W_0\boldsymbol{g}_0 + \boldsymbol{\varepsilon}$, where $\boldsymbol{y}$ is the same as before, $\boldsymbol{g}_0$ is an $m \times 1$

383    unobservable random vector containing allele substitution effects of each marker, $(W_0)_{n\times m}$ is the

384    random observable design matrix which is of the form $(W_1' \vdots \cdots \vdots W_S')'$ when ordering data by

385    subpopulation, and $\boldsymbol{\varepsilon}$ is a random vector of residuals. The priors for $\boldsymbol{g}_0$ are simply univariate versions

386    of the priors used for $\boldsymbol{g}$. Thus, $\boldsymbol{g}_0|G^D \sim \pi(\cdot\ |G^D), G^D = Diag\left(\sigma_{g1}^2, \dots, \sigma_{gm}^2\right), \sigma_{gj}^2 \overset{iid}{\sim} IG\left(\frac{a}{2}, \frac{b}{2}\right)$, (for the

387     homogeneous marker effect variance model $\sigma_{g1}^2 = \cdots = \sigma_{gm}^2 = \sigma_g^2$) and the residual variance $\sigma^2$ is

388     given an $IG\left(\frac{\tau^2}{2}, \frac{v}{2}\right)$ prior as before. In addition, $\boldsymbol{p} = (p_1, p_2, \ldots, p_m)$ is a vector of overall reference

389     allele frequencies, $W_0|\boldsymbol{p} \sim \pi(W_0|\boldsymbol{p})$ is a simplified version of $\pi(W|P^*)$ (shown later), and the prior

390     for $\boldsymbol{p}$ is $p_j \overset{iid}{\sim} Beta(\alpha, \beta), j = 1, 2, \ldots, m.$

391     The Deviance Information Criterion (DIC; Spiegelhalter et al., 2002) combines a measure of

392     goodness of fit based on the posterior distribution and a penalty for model complexity, and despite

393     some criticism it has been used in different areas to perform model comparison (Gelman et al., 2014;

394     Spiegelhalter et al., 2014). It has the following form:

$$DIC = -2\log f\left(\boldsymbol{Data}|\widehat{\boldsymbol{\theta}}_B\right) + 2p_{DIC}$$

395     where $p_{DIC} = 2\left(\log f\left(\boldsymbol{Data}|\widehat{\boldsymbol{\theta}}_B\right) - E_{\boldsymbol{\theta}|\boldsymbol{Data}}[\log f(\boldsymbol{Data}|\boldsymbol{\theta})]\right), \widehat{\boldsymbol{\theta}}_B = E[\boldsymbol{\theta}|\boldsymbol{y}]$ is the posterior mean of

396     the unknown parameters. The first component of $DIC$ is a measure of model adequacy, whereas the

397     second one is the effective number of parameters which is a penalty for increasing model complexity

398     (Spiegelhalter et al., 2002). Models with a smaller DIC are preferred. Recall that for any of our

399     models the likelihood has two components: $f(\boldsymbol{y}, W^\sigma|W^N, \boldsymbol{g}, R, P^*) = f(\boldsymbol{y}|W, \boldsymbol{g}, R)f(W^\sigma|W^N, P^*)$

400     that were denoted as the $\boldsymbol{y}$ component and the $W$ component. Thus, the general form of the DIC is:

$$DIC = -2\log f\left(\boldsymbol{y}|W^\sigma, \widehat{W}_B^N, \widehat{\boldsymbol{g}}_B, \widehat{R}_B\right) + 2p_{DIC-y} - 2\log f\left(W^\sigma|\widehat{W}_B^N, \widehat{P}_B^*\right) + 2p_{DIC-W}$$

$$:= DIC_{\boldsymbol{y}} + DIC_W$$

401     where $p_{DIC-y} = 2\left(\log f\left(\boldsymbol{y}|W^\sigma, \widehat{W}_B^N, \widehat{\boldsymbol{g}}_B, \widehat{R}_B\right) - E_{W^N, \boldsymbol{g}, R, P^*|\boldsymbol{y}, W^\sigma}[\log f(\boldsymbol{y}|W, \boldsymbol{g}, R)]\right)$ and $p_{DIC-W} =$

402     $2\left(f\left(W^\sigma|\widehat{W}_B^N, \widehat{P}_B^*\right) - E_{W^N, P^*|\boldsymbol{y}, W^\sigma}[f(W^\sigma|W^N, P^*)]\right)$. Thus, as the likelihood, the DIC can be

403     decomposed into a $\boldsymbol{y}$ component $DIC_{\boldsymbol{y}}$ and a $W$ component $DIC_W$.

404     *2.4 Parameter inference via MCMC*

405  In this section, some issues about MCMC algorithms to carry out inference are briefly discussed.

406  Notice that when $W$ is fully observed, the fact that there are no missing genotypes implies that

407  posterior sampling for the (hyper) parameters of the $W$ component of the likelihood and the (hyper)

408  parameters of the $\boldsymbol{y}$ component can be performed separately. The full conditionals of $\boldsymbol{g}, G, \sigma^2, g_0$, and

409  $\sigma_g^2$ are known; therefore, samples from the joint posterior can be obtained using a Gibbs sampler

410  (Casella and George, 1992) while samples from the posterior distribution of allelic frequencies can

411  be obtained using a Metropolis-Hastings algorithm. Specifically, independent Metropolis algorithms

412  are considered here. For the scenario of $\boldsymbol{r}$ known, the new samples can be generated in two steps:

413  firstly a Dirichlet vector is sampled, and secondly its elements are scaled with the appropriate

414  elements of $\boldsymbol{r}$. Alternatively, uniform$(0, r_l)$ distributions can be used as proposal, which simplifies

415  computations. With such proposal, given the current state of the chain denoted as $P^t$, the acceptance

416  probability of the new sample $P_+^t$ is min $\left\{\frac{\pi(P_+^t|W)}{\pi(P^t|W)}, 1\right\}$. For null models, the posterior distribution of

417  $\boldsymbol{p}_0$ is the product of $m$ Beta$(p_j; n^{B_j} + \alpha, n^{A_j} + \beta)$ distributions, $j = 1,2 \dots, m$. Hence, direct

418  sampling can be implemented if needed and the functional form of the posterior mean is known.

419  When $\boldsymbol{r}$ is unknown, the candidate to sample from the posterior of $(\boldsymbol{p}_j, \boldsymbol{q}_j), j = 1,2, \dots, m$, could be a

420  Dirichlet distribution.

421  On the other hand, when matrix $W$ is partially observed a Metropolis-within-Gibbs strategy (Robert

422  and Casella, 2010) can be used to sample from the joint posterior. This strategy is useful due to the

423  fact that nor $\pi(W^N|Else)$ neither $\pi(P^*|Else)$ are standard distributions and the existence of the

424  parameter $W^N$ does not allow to carry out separate sampling algorithms as before because this is a

425  parameter of both components of the likelihood. Accordingly, there are two Metropolis steps in the

426  algorithm to sample from the posterior of the full models. The first one is used to obtain samples

427  from $\pi(W^N|Else)$. A good proposal is $\pi(W^N|W^o, P^*)$ because obtaining direct samples from this

428  distribution via the inverse transform method for discrete random variables (Robert and Casella,

429    2010) is straightforward. The functional form of $\pi(W^N|W^\sigma, P^*)$ is derived from first principles as

430    explained in 2.3.1. Thus, given the current state of the chain $W^{N_t}$, the acceptance probability of a

431    new sample $W_+^{N_t}$ is: $\min\left\{\frac{\pi(W_+^{N_t}|Else)\pi(W_+^{N_t}|W^\sigma, P^*)}{\pi(W^{N_t}|Else)\pi(W^{N_t}|W^\sigma, P^*)}, 1\right\}$. This applies to both situations: $\boldsymbol{r}$ known and $\boldsymbol{r}$

432    unknown.  The second Metropolis step is used to draw samples from $\pi(P|Else)$ for $\boldsymbol{r}$ known or

433    $\pi(P, Q|Else)$ for $\boldsymbol{r}$ unknown.  The proposals mentioned for the non-missing genotypes scenario also

434    work here. For the null models, it turns out that $\forall\, j = 1,2,\ldots,m, \pi(p_j|Else)$ is a known distribution,

435    it is a Beta($n^{Bj} + \alpha, n^{Aj} + \beta$) and consequently only one Metropolis step is needed because direct

436    sampling from the full conditional distribution of $\boldsymbol{p}_0$ is feasible.  Notice that this full conditional

437    distribution is the posterior distribution of $\boldsymbol{p}_0$ when matrix $W$ is completely observed.

438    *2.5 Simulation study*

439    In order to provide an example of the implementation of some of the proposed models and the

440    computation of some criteria to compare their performance, two simulated datasets were used.

441    Simulation of these datasets involved two main steps: Simulation of genotypes (QTL and SNP), and

442    simulation of QTL effects and noise. The phenotypes were simulated as the sum of additive genetic

443    effects (sum of QTL allele content times the allele effect) and noise. Datasets were simulated using

444    the software QMSim (Sargolzaei and Schenkel, 2013). In both cases, a historical population was

445    simulated by creating 1000 generations of random mating using a forward-in-time approach in order

446    to reach mutation-drift equilibrium and to create linkage disequilibrium (Sargolzaei and Schenkel,

447    2013). The historical population size in each generation was 1000 with 500 males and 500 females.

448    Then, subpopulations were created from individuals pertaining to the historical population under

449    different selection pressures and criteria, and different mating systems (Table 1).

450    Phenotypes were simulated with different number of QTL controlling the trait and different

451    heritabilities. Furthermore, the population structure also differed because the criteria to simulate the

452    subpopulations were different for each trait. Briefly, dataset 1 involved three subpopulations with

453    different number of generations, migration was allowed and the heritability of the trait was high.

454    Dataset 2 comprised two subpopulations with only two generations, no migration and the heritability

455    of the trait was low (Table 1). For further details concerning the simulation see appendix B.

456    Given that this paper is focused on proposing and explaining a set of across population genome-wide

457    prediction models and not with their large scale implementation, the number of simulated SNP and

458    sample size were low in order to avoid computational issues (Table 1). Phenotype 1 illustrates the

459    situation in which the number of markers is equal to the number of QTL affecting the trait, while for

460    phenotype 2 the number of markers is larger than the number of QTL controlling the trait. These

461    contrasting simulation schemes, different selection pressures and criteria, mating designs and number

462    of generations were used to mimic real life situations where different subpopulations have different

463    backgrounds. These simulated datasets were used to carry out analyses using the following models:

464    Homogeneous and heterogeneous marker effect covariance matrices with homoscedastic residuals

465    and their null versions. Only models with homoscedastic residuals were used to analyze these

466    datasets because simulations did not consider heteroscedastic residuals.

467     The analyses performed involved implementation of MCMC algorithms explained in section 2.4, the

468    computation of DIC and the computation of the following quantities measuring predictive

469    performance and accuracy: the squared correlation between predicted breeding values and

470    phenotypes in the testing populations, hereinafter called predictive ability, and squared correlations

471    between true and predicted breeding values computed in the testing populations (accuracy). Because

472    true breeding values were available for the complete populations, squared correlations between true

473    and predicted breeding values in the training populations were also computed.

474    For dataset 1, the training population was comprised of generations 0 to 2 of subpopulation 1, 0 to 5

475    from subpopulation 2 and generation 0 of subpopulation 3, while the testing population included

476    generation 3 of subpopulation one, generation 6 of subpopulation 2 and generation 1 of

477    subpopulation 3. For dataset 2, the training population was composed of generations 0 and 1 of

478    subpopulations 1 and 2 and the testing dataset contained generation 2 of subpopulations 1 and 2.

479    In dataset 2, the full genotypes of three individuals (one founder from each subpopulation and a non-

480    founder from subpopulation 1) were not included in the analysis in order to simulate the case of

481    missing genotypes.

482    It was assumed that $r = \left(\frac{1}{S}, \ldots, \frac{1}{S}\right)$. In an initial analysis, a scaled Dirichlet distribution was used as

483    proposal to draw samples from $\pi(P|Else)$, but the behavior of the chains was not satisfactory

484    because the acceptance rate was too low (results not shown). Consequently the product of $S$

485    independent uniform $\left(0, \frac{1}{S}\right)$ distributions was used as proposal. For each dataset, 20.000 iterations

486    were run; the first 10.000 were considered burn-ins. An in-house R script (R Core Team, 2015) was

487    created to carry out the analyses which were performed using the University of Florida's high

488    performance computing cluster.

489    **3. Results**

490    *3.1 Simulated populations*

491    Tables 1 and 2 show features corresponding to characteristics of the simulated genomes and

492    populations.

493    **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[1] | 1: 28M, 180F, Phen/L<br>2: 20M, 90F, Phen/H<br>3: 50M, 500F, Rnd | 1: 5M, 25F, Rnd<br>2: 20M, 50F, Phen/H |

| Number of generations, mating system and selection criteria[2] | 1:3,0.8,0.4, As1/Phen, Phen/L<br>2: 6, 0.7, 0.1, As2/Phen, Phen/H<br>3:3, 0.7, 0.2, Rnd, Rnd | 1: 2, 1, 0.9, Rnd, Rnd<br>2: 2, 0.9, 0.3, Rnd, Phen/H |

[1]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).

[2]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 1).

**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)<br>S2: 25.71 (9.60)<br>S3: 0.26 (9.91) | S1:-0.5959 (9.3616)<br>S2:8.9253 (11.9571) |

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 4 shows DIC values for dataset 1, Table 5 DIC values for dataset 2 and Table 6 shows predictive abilities and accuracies in both datasets. For Tables 4 to 6, the following is the meaning of abbreviations for the different models fitted to datasets 1 and 2: $M_{1G}$= full model with Multivariate Gaussian prior and homogeneous marker effect covariance matrices, $M_{1G}^*$= full model with Multivariate Gaussian prior and heterogeneous marker effect covariance matrices. Recall that all models assumed homoscedastic

514    residuals. The remaining models with subindex 1 replaced by 0 correspond to null versions of the

515    corresponding full models.

516    **Table 4 $y$** component and total DIC for dataset 1

| Model | $y$ component of DIC | Total DIC |
|:---:|:---:|:---:|
| $M_{1G}$ | 33702.55 | 4751373.55 |
| $M_{1G}^*$ | 11599.05 | 4729270.05 |
| $M_{0G}$ | 15396.32 | 6604501.32 |
| $M_{0G}^*$ | 13008.42 | 6602113.42 |

517    Thus, in dataset 1, according to the $y$ component of DIC, for the models with homogeneous marker

518    effect covariance matrices (variances) the null model performed better, while for models with

519    heterogeneous covariance matrices (variances) according to this criterion the full model should be

520    preferred over its null version. When considering the whole likelihood to compute the DIC, the two

521    full models had smaller DIC. Additionally, the model with the smallest DIC, and therefore the "best"

522    one under this criterion was model $M_{1G}^*$.

523    **Table 5 $y$** component, $W$ component and total DIC for dataset 2

| Model | $y$ component of DIC | $W$ component of DIC | Total DIC |
|:---:|:---:|:---:|:---:|
| $M_{1G}$ | 1314.0 | 38367.4 | 39681.4 |
| $M_{1G}^*$ | 1328.8 | 38356.4 | 39684.2 |
| $M_{0G}$ | 1365.6 | 38180.3 | 39545.9 |
| $M_{0G}^*$ | 1370.1 | 38179.0 | 39549.1 |

524    In this dataset the two components of the DIC values and therefore DIC values were similar for all

525    models. The $y$ components of DIC were smaller for the full models.  Conversely, the $W$ components

526    were smaller for null models as well as total DIC values.

527    **Table 6** 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 |

528  In dataset 1, according to predictive abilities, the model with the best performance was model

529  $M_{0G}^*$ while model $M_{1G}$ had the worst performance. The squared Pearson correlations between true

530  and predicted breeding values in testing dataset 1 suggested that the performance of these models

531  followed a trend similar to that indicated by predictive abilities. In training dataset 1, model

532  $M_{1G}^*$ yielded the highest accuracy and model $M_{1G}$ had the smallest accuracy.

533  Predictive abilities and accuracies in the testing sets were extremely low for dataset 2. Accuracies in

534  training set were higher than those obtained in the testing set; however, they were still low. There

535  were not substantial differences between these squared correlations. Predictive abilities were higher

536  for the full models, while accuracies in testing and training sets were higher for the null models.

537  **4. Discussion**

538  4.1 General features of the models

539  A group of hierarchical Bayesian linear regression models to carry out simultaneous genome-wide

540  prediction in several subpopulations accounting for randomness of genotypes was presented. The

541  proposed models differed in the prior distribution assigned to the marker effects and on the

542  assumptions made about residual variances (homogeneous or heterogeneous across subpopulations).

543  The priors for the marker effects were multivariate (univariate) Gaussian and allowed homogeneous

544  or heterogeneous covariance matrices (or variances).

545  The differences between these models and other regression models currently used in across

546  population genome-wide prediction are: 1) subpopulation-specific effects for each marker are

547  considered and their covariance matrices are modeled explicitly, and 2) genotypes are treated as

548  random variables with a distribution that depends on allelic frequencies as well as on pedigree

549  information. The second feature makes these models different from all other genome-wide prediction

27

550  models. The distribution of genotypes combines pedigree and genomic information that are not used

551  when randomness of $W$ is ignored. It allows accounting for heterogeneity and correlations of allelic

552  frequencies of the same marker across subpopulations and including individuals with phenotypes and

553  missing genotypes in various loci without carrying out a previous imputation. This is possible

554  because the non-observed part of $W$, denoted as $W^N$, is treated as a parameter and therefore

555  imputation is automatically performed. Another advantage is that the use of a Bayesian approach

556  automatically takes into account uncertainty about the imputed genotypes.

557  Although most of the paper has been devoted to the models allowing subpopulation-specific effects

558  for each marker (the full models), their univariate versions (the null models) are also contributions of

559  this study. These also allow including individuals with missing genotypes in some or all marker loci

560  without need of external imputation and take into account randomness in genotypes. Therefore, these

561  models could also be used either in single population analyses or to conduct across population

562  genome-wide prediction pooling the data as has been done in previous studies (de Roos et al., 2009;

563  Lund et al., 2011; van den Berg et al., 2015; Wientjes et al., 2015) and was also done here.

564  Doing a joint analysis has the advantage that the number of phenotypes increases, but in our full

565  models the number of location parameters is also incremented because each marker is allowed to

566  have subpopulation-specific effects; moreover, the number of covariance parameters also increases.

567  The gain in accuracy is achieved when factors such as different QTL effects across subpopulations,

568  differences in linkage phase between QTL and markers, and differences in allelic frequencies and LD

569  patterns make marker effects change substantially from one subpopulation to another. Consequently,

570  the performance of these models may have considerable variation from one dataset to another.

571  The diagonal blocks of $G$ were assumed to be non-structured. A way reduce dimensionality of the

572  parameter space is to assume certain structure of $G$. For example, it can be assumed that all

573  covariances and variances are the same, thus, only two parameters per block have to be estimated.

574      The conditional independence property used to derive $\pi(W|P^*)$ implies that allelic frequencies are

575      estimated in the set of oldest individuals with phenotypes. Here, this set of individuals was referred

576      to as the base population and individuals pertaining to it were referred to as founders. This was done

577      for pragmatic purposes. However, truncating the pedigree by ignoring individuals without phenotypic

578      records created a group of individuals that may not be the actual base population which is defined as

579      that comprised by ancestors with unknown parents (Henderson, 1974; Kennedy et al., 1988).

580      Conversely, in other cases phenotypic records from this population may be available; thus, estimates

581      of allelic frequencies in the true base population can be obtained. Here, it was further assumed that

582      founders were unrelated which is likely to be false in many situations. However, this assumption has

583      been made in conventional models used to do genetic analysis (Henderson, 1974; Kennedy et al.,

584      1988) because pedigrees are not always completely known. Consequently, what is called the base

585      population is not always the true one. Nevertheless, this assumption seems to be reasonable after so

586      many years of successful artificial selection in animals and plants based on predicted breeding values

587      obtained from these models (Hill, 2014; Gianola and Rosa, 2015).

588      As discussed in section 2.1.1, the pmf $\pi(W|P^*)$ could be derived ignoring pedigree information.

589      Then, this pmf could be found as the product of all $\pi\left(w_{ij}^l \big| p_{lj}^*\right)$ or the product of binomial

590      distributions for gene content (i.e., the number of copies of the reference allele at each locus) across

591      loci and individuals with each binomial distribution depending on the corresponding allelic

592      frequencies. Notice that this requires reparametrizing the mapping of genotypes, that is, instead of

593      having $\{-1,0,1\}$ as possible values of an entry of $W$, values would be $\{0,1,2\}$. In this case, all

594      individuals in the population would be used to estimate allelic frequencies instead of using

595      information from a base population. If pedigree information is available, it can be easily incorporated

596      into the derivation of $\pi(W|P^*)$ as was shown here and the resulting pmf is not very difficult to

597      evaluate. Furthermore, as mentioned before, direct sampling from this pmf can be done via the

598    inverse transform method for discrete random variables. Notwithstanding, in scenarios where

599    pedigree information is very scarce or not reliable, adding the assumption of independence among

600    individual genotypes and using binomial distributions for the gene content of each individual at each

601    marker locus is an option to model the distribution of matrix $W$ which would induce a joint pmf

602    similar to those presented in Gianola et al., (2010) and Martínez et al. (2015).

603    If some individuals with phenotypes have only one known parent, the pmf of their genotypes

604    conditioned on this parent and allelic frequencies can be defined in a similar way as was done in

605    Table A.1 for the case of a fully known pedigree (see Appendix C). In this situation, *Remark 1* does

606    not hold and the functional form of $\pi(W|P^*)$ changes which implies that $\pi(W|Else)$ changes as

607    well.

608    Regarding assumptions about the distribution of allelic frequencies, our models allow for correlations

609    between them. To do that, priors based on a Dirichlet distribution were used. Using these priors

610    require allelic frequencies to be expressed on a complete population basis. This setting brings

611    parameter $\boldsymbol{r}$ into the picture. The algebra associated with this parameter is clear and straightforward,

612    but its interpretation may be fuzzy. From an algebraic standpoint, these parameters are upper

613    boundaries posed over allelic frequencies to force them to be in the support of the prior distribution,

614    thus they can be seen as analytic instruments. Nevertheless, their meaning from the population

615    genetics standpoint is not very clear. Perhaps, the easier interpretation when assuming $r_{1l} = \cdots =$

616    $r_{ml} = r_l$, is that $r_l$ is the relative frequency or weight of the $l^{th}$ subpopulation. However, making

617    claims about the biological interpretation of this set of parameters is beyond the scope of this study.

618    From a statistical viewpoint, two approaches were proposed. The first one assumed that $\boldsymbol{r}$ was known

619    (truly known or set to some *ad hoc* value) and $r_{1l} = \cdots = r_{ml} = r_l$. In the examples used here all

620    subpopulations were given the same weight, that is, $r_l = 1/S, \forall\, l = 1,2 \ldots, S$, a pragmatic decision

621    that has been used in other studies, e.g., Gianola et al. (2010). In this scenario, for all $j$, $\boldsymbol{p}_j$ is modeled

622  as a scaled Dirichlet vector which allows non-null covariances between its elements. The second

623  approach assumed that $r$ was unknown and $\{r_{lj}\}$ varied across marker loci. For each locus the prior

624  was a Dirichlet over allelic frequencies of both alleles in all subpopulations and it permitted

625  obtaining posterior samples of allelic frequencies and $r$. Under the assumption of independence of

626  allelic frequencies, independent priors could be assigned to each marker (e.g., Uniform$(0,r_l)$) and the

627  validity of this assumption could be tested using criteria as Bayes factors or DIC. If data are pooled

628  and structure is ignored (as done in the null models) the full conditional pdf $\pi(\boldsymbol{p}_0|Else)$ is known

629  and therefore direct sampling can be implemented when matrix $W$ is not completely observed. On

630  the other hand, when it is completely observed the posterior of $\boldsymbol{p}_0$ is known and there is no need of

631  sampling to obtain point estimators. The reason for the full conditional of $\boldsymbol{p}_0$ being a known

632  distribution but not its posterior in the presence of missing genotypes is that $W^N$ is an extra

633  parameter in the model and obtaining the marginal posterior of $\boldsymbol{p}_0$ implies marginalization of

634  $\pi(W^N,\boldsymbol{p}_0|W^o)$ over $W^N$ which induces a non-standard pmf.

635  The derivation of the pmf $\pi(W|P^*)$ and $\pi(W_0|\boldsymbol{p}_0)$ not only allow inferences concerning the marker

636  allelic frequencies in the base population, but also allow predictions for non-genotyped or partially

637  genotyped animals without performing a previous imputation. This is likely to increase accuracy of

638  genome-wide predictions because it allows incorporating more phenotypic records. Imputed missing

639  genotypes can be obtained using posterior means or medians of $W^N$. However, these outputs have to

640  be viewed as a byproduct because these models were not intended to perform imputation. The

641  imputation of missing genotypes is an underlying process in the prediction of genotypic values of

642  individuals with missing genotypes. Notwithstanding, because samples from the posterior of $W^N$ are

643  available and computation of imputed genotypes is simple, there could be interest in using this output

644  of the model and in such case the accuracy of the imputation would also be of interest. Hence,

645  although imputation was not a main objective of our models, it is worth making a brief comment on

646    it. Though an assessment of imputation accuracy is a matter for further research, two statements can

647    be made about the imputation process in our models. Firstly, one advantage of the models developed

648    here is that they automatically take into account the uncertainty of imputation (as a consequence of

649    using a Bayesian approach). Conversely, in the standard approach where genotype imputation is the

650    first step and then a random linear regression model is fitted using these imputed values as if they

651    were observations, uncertainty is not taken into account. Secondly, a disadvantage of our models is

652    that they do not incorporate LD information when imputing missing genotypes, a source of

653    information that is used by some of the current imputation methods (Li et al., 2009). Here, pedigree

654    information, phenotypes and allelic frequencies are used for imputation. Thus, benchmarking of the

655    procedure developed here with current and well-accepted procedures is material for future studies.

656    Furthermore, another question that can be addressed in future research is if improving this imputation

657    as discussed later in section 4.3 has a significant impact on the predictive performance of the models.

658    As mentioned before, the regression models used in genome-wide prediction treat genotypes as fixed

659    and their effects as random while in the classical quantitative genetics theory genotypes are treated as

660    random and allelic substitution effects as fixed. The set of models developed here are something in

661    between because genotypes are treated as random variables as in classical quantitative genetics, and

662    marker effects are considered random as well like in the standard regression models used in genome-

663    wide prediction. de los Campos et al. (2015b) presented an excellent discussion on the connections

664    between the heritability and the so-called genomic heritability obtained with linear regression

665    models. They show why caution has to be exercised when interpreting the parameters obtained using

666    genomic information due to the fact that sometimes the connection between parameters as the

667    additive genetic variance and the genomic variance are not straightforward. Similarly, Gianola et al.

668    (2015) discussed the fact that connections between genomic correlations and additive genetic

669    correlations are ambiguous. So far, the Bayesian models proposed in this paper are intended to

670    predict breeding values, phenotypes, and to estimate allelic frequencies in a base population using

671 genomic information and no claim is made about the properties of covariance parameters obtained

672 from them.

673 The discussion above is relevant because the regression variables are not based on genes, but proxies

674 for the causal variants affecting the phenotypes of interest. However, taking into account these

675 limitations and the high degree of caution needed when interpreting parameters obtained from

676 models using molecular markers, some parameters such as the fraction of additive genetic variance

677 explained by the markers are of interest and our models could be used to estimate these quantities.

678 The family of models developed here could be applied or adapted to different situations. In the

679 simulation, the case of individuals coming from a common founder population pertaining to

680 subpopulations with different selection criteria and mating systems was considered. Other situations

681 in which this set of models could be useful are: 1) simultaneous evaluation of individuals from

682 different breeds or lines, 2) individuals from the same breed or line performing under different

683 environmental conditions (e.g., different geographic regions, production systems, etc.), 3) a

684 combination of numerals 1 and 2, 4) simultaneous evaluation of several correlated traits. In this last

685 case, if all individuals have records for all phenotypes, the design matrix satisfies $W = I_S \otimes W_+$,

686 where $W_+$ is the matrix of dimension $n \times m$ containing genotypes of $n$ individuals at $m$ marker loci.

687 In this case the model is being adapted to handle correlations between the effects of a given marker

688 locus for different traits in a single population. Consequently, for a given choice of prior and

689 assumption about residuals (heteroscedastic or homoscedastic) the model involves the corresponding

690 hierarchical structure except for the pmf of $W$ conditional on the allelic frequencies and pedigree

691 which is $\pi(W_+|\boldsymbol{p}_0^*)$ instead of $\pi(W|P^*)$. Recent studies have developed Bayesian multiple-trait

692 genome-wide regression models and have shown that predictions from them are more accurate than

693 those coming from genomic univariate models (Jia and Jannink, 2012). The hierarchical Bayesian

694 multivariate genome-wide prediction models proposed by Jia and Jannink (2012) have similar

695  components to the models presented here such as the priors for $\boldsymbol{g}$, but they do not account for

696  randomness of genotypes. Another step to accommodate our models for multiple-trait prediction is to

697  allow correlated residuals, that is, a non-diagonal matrix $R$. In this case, an inverse Wishart prior can

698  be assigned instead of the inverse gamma prior used here.

699  4.2 Simulation results

700  As stated in section 2.4, the aim of this limited simulation was to provide an illustration of the

701  implementation of models and methods developed in this study. Thus, results are not conclusive and

702  further research involving analyses based on more elaborate simulations as well as real datasets to

703  have a better evaluation of the performance of this family of models is needed. Nevertheless, some

704  insights and comments derived from the analyses of these two datasets can be discussed.

705  The correlation between phenotypes and predicted breeding values (or its square) is one of the most

706  widely used measurements to compare genome-wide prediction models, it is associated with the

707  response to selection and it is easy to compute. On the other hand, as mentioned previously, the DIC

708  combines measures of model adequacy and complexity (Spiegelhalter et al., 2002).

709  For dataset 1, the squared correlation between phenotypes and predicted breeding values (the

710  predictive ability) did not show an advantage in predictive capability of models taking into account

711  the population structure, i.e., the existence of the subpopulations (Table 6). While measures based on

712  squared correlations did not provide conclusive evidence in favor of the full models, the DIC favored

713  the full models.

714  As expected, the predictive ability and the other correlations were much smaller in dataset 2 due to

715  the lower heritability of the trait. Although all predictive abilities were low, according to this

716  criterion the performance of the full models was slightly better. Accuracies of predicted breeding

717  values suggested a tiny superiority of null models. The two subpopulations simulated in this dataset

718  diverged by just two generations which could cause only small differences in allelic frequencies, this

719    scenario clearly favors the null models. Accordingly, the DIC component coming from genotypes

720    was slightly better (smaller) for null models as opposed to the case of dataset 1. The total DIC gave

721    evidence in favor of null models. Among predictive ability, accuracy and DIC, accuracy and DIC

722    favored the null models, but the values were very close. The performance of the fitted models was

723    more similar in this dataset than in dataset 1.

724    In our small simulations, when subpopulations diverged by several generations, migration was

725    allowed and heritabilities were high (dataset 1), full models had better performance in terms of DIC.

726    Conversely, when populations diverged by only a few generations, there was no migration, and

727    heritabilities were low (dataset 2) null models tended to perform better according to this criterion.

728    However, the differences were small. On the other hand, predictive abilities showed a different

729    pattern. In dataset 1 this criterion was higher for null models while in dataset 2 it was smaller for null

730    models. Another feature shown by these simulations was the high variability in model performance

731    that may exist among populations. In dataset 1, according to all criteria except the $W$ component of

732    DIC, the performance of model $M_{1G}$ tended to be remarkably poorer while this was not the case in

733    dataset 2.

734    Other authors have found modest or null increments in predictive performance of models allowing

735    heterogeneous marker effects across subpopulations compared to pooling data and analyzing the

736    complete population as a single one (Olson et al., 2012; Makgahlela et al.,2013; de los Campos et al.,

737    2015a). All the aforementioned studies used real data from plants and animals. Working with three

738    plant populations and using a model very similar to those proposed here, Lehermeier et al. (2015)

739    found cases in which the strategy of pooling data and ignoring structure performed better and other

740    cases where multivariate models yielded better predictive performance. These authors found that in

741    highly differentiated populations within group and multivariate analyses performed better while the

742      converse occurred in closely related subpopulations with small sample sizes. Roughly speaking,

743      these results are in agreement with the results found in this study.

744      Using predictive ability, Lund et al. (2011) found a higher accuracy of predicted additive breeding

745      values when pooling the data compared with individual analyses. Similar results were found by de

746      Roos et al. (2009) when heritability was low, divergence of populations was small (small number of

747      generations) and marker density was high (more persistent phase), and by Wientjes et al. (2015)

748      when the QTL effects did not change across subpopulations. Pooling data and ignoring the

749      population structure corresponds to the null models defined in this study, except that models

750      considered by the authors just cited did not account for randomness of genotypes. In our simulation,

751      individual analyses were not considered. Sample size is one the factors affecting the accuracy of

752      genome-wide predictions (Meuwissen et al., 2001; Goddard 2009, Zhong et al., 2009). Presumably it

753      was one of the leading factors causing the results found by Lund et al. (2011). In addition, the

754      Holstein breed is highly inbred and there were several individuals connecting the different

755      populations; this probably made them similar. On the other hand, the studies of de Roos et al. (2009)

756      and Wientjes et al. (2015) used simulated data and explored different scenarios. Both studies found

757      situations in which pooling data was not advantageous.

758      4.3 Refinements and extensions

759      In this section, some comments regarding possible extensions and refinements of different aspects of

760      the family of models presented in the study are briefly discussed.

761      In the derivation of the joint pmf of $W$ conditional on $P^*$ and pedigree information, row-wise

762      dependence due to kinship was taken into account by using pedigree information to accommodate

763      relationships among genotypes of related individuals. This task was highly simplified due to the

764      conditional independence argument that permitted to find a simpler decomposition of the joint pmf

765      and therefore, a simpler algebraic expression. However, the possible existence of column-wise

766 dependence due to LD was ignored here in order to make the problem more tractable from the

767 mathematical point of view. This is an assumption frequently used in theoretical studies in

768 quantitative genetics and it is well-accepted at least in studies concerned with first approximations to

769 a given problem. For example, Gianola et al. (2009) treated a series of theoretical aspects of some of

770 the Bayesian regression models used in genome-wide prediction using the assumption of linkage

771 equilibrium which implies the mutual independence of the columns of $W$ used here (they also

772 developed some results accounting for LD in the Appendix). Most of the models currently used in

773 genome-wide prediction are also based on this assumption, few approximations to deal with

774 consequences of LD have been proposed (Gianola et al., 2003; Yang and Tempelman, 2012), but

775 these have not yet been adopted in routine genetic evaluations. Their models do not consider

776 randomness in the genotypes; thus, a consequence of considering LD in these models is the need to

777 account for covariances between marker effects at different loci. Consequently, a refinement of our

778 family of models in this regard, would be to accommodate LD, which can be performed at two

779 levels: 1) account for correlations among columns of $W$, and 2) use a non-block-diagonal $G$ matrix.

780 A potential consequence of accounting for non-independence of the columns of $W$ could be the

781 reduction in the cardinality of $\mathcal{G}$ that is induced by the fact that the number of possible values of a

782 column of $W$ depends on the values at one or more different columns (as it happened with rows).

783 Another assumption made here was the absence of mutations which caused that when conditioning

784 on the genotypes of the parents of an individual, the probabilities of its genotype taking a given value

785 were completely defined by the parental genotypes, making this random variable conditionally

786 independent of allelic frequencies. Thus, another refinement in $\pi(W|P^*)$ would be to account for

787 mutation. Therefore, the derivation of $\pi(W|P^*)$ to accommodate dependence between columns of $W$

788 and mutation, and the impact of this refinement on predictive performance and the accuracy of

789 imputed genotypes (if it is of interest) pose a problem for further research.

37

790    If relationships among founders (as defined in this paper) were to be taken into account, from the

791    theoretical point of view it is not hard to visualize how to do it. For the sake of simplicity, the case of

792    two individuals and one locus is considered; consequently, the sub-index associated with locus is

793    omitted. Let $W_1, W_2$ be the genotypes of individuals 1 and 2, and $W_{\mathcal{C}}$ the genotypes of the set of

794    relevant common ancestors. Suppose that 1 is not a parent of 2. Then:

795    $\pi(W_1, W_2 | P^*) = \sum_{\mathcal{G}^{\mathcal{C}}} \pi(W_1, W_2 | W_{\mathcal{C}}, P^*) \pi(W_{\mathcal{C}} | P^*) = \sum_{\mathcal{G}^{\mathcal{C}}} \pi(W_1 | W_{\mathcal{C}}, P^*) \pi(W_2 | W_{\mathcal{C}}, P^*) \pi(W_{\mathcal{C}} | P^*),$

796    where $\mathcal{G}^{\mathcal{C}}$ is the set of possible values that the set of genotypes of relevant common ancestors can

797    take according to the pedigree (as explained in section 2.1.1) and the second equality follows from

798    the conditional independence of the genotypes of individuals 1 and 2 given the common ancestors

799    and allelic frequencies. By relevant common ancestors it is meant that the genotypes of these

800    ancestors provide information about the genotypes of 1 and 2 when conditioning on the full set of

801    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-

802    index $c$ means complement with respect to $\mathcal{D}$) and $\pi(W_1, W_2 | W_{\mathcal{D}}, P^*) = \pi(W_1, W_2 | W_{\mathcal{C}}, P^*)$. Notice

803    that unless individuals 1 and 2 are full sibs, their conditional pmf given the relevant common

804    ancestors depends on $P^*$. Of course, it makes $\pi(W | P^*)$ a more complex expression and reduces the

805    cardinality of $\mathcal{G}$. See Appendix D for a toy example of $\pi(W_1, W_2 | W_{\mathcal{C}}, P^*)$ when 1 and 2 are half sibs.

806    Although the problem is tractable from the theoretical standpoint, it may be difficult to compute

807    these values especially with complex pedigrees where the set of common ancestors may be large

808    such as those found in animal and plant populations. The example in Appendix D shows that even in

809    a simple case, computation of $\pi(W_1, W_2 | P^*)$ is involved.

810    **5. Conclusions**

811    The main contribution of this paper is the theoretical development of a set of models for across

812    population genome-wide prediction incorporating marker genotypes not only as explanatory

813    variables of regression models, but also as realizations of random variables providing information

814     about allelic frequencies and missing genotypes. Although models were intended for across

815     population analysis, they can also be applied in single population studies and adapted for multiple-

816     trait prediction.

817     Theoretical and computational issues along with possible applications as well as some extensions and

818     refinements of these models pose several problems for future research. Our models treat both

819     genotypes and marker allelic substitution effects as random; therefore, they combine features from

820     classical quantitative genetics theory and traditional genome-wide prediction models.

821     Some features of the models developed in this study make them promising for genome-wide

822     prediction. Among these, the ability to include phenotypes from individuals with missing genotypes

823     at some or all loci without the need of previous imputation and accounting for uncertainty about

824     imputed genotypes as well as heterogeneity of allelic frequencies across subpopulations are perhaps

825     the most appealing. Further research to assess their performance and also to compare them with other

826     models used in genome-wide prediction is needed.

827     **Author Contributions**

828     C.A. Martínez developed modeling strategies, carried out the derivations, wrote the R scripts,

829     designed and made the simulations and wrote the paper. K. Khare advised modeling strategies,

830     reviewed, corrected and discussed the derivations and the statistical aspects of the paper. A. Banerjee

831     advised modeling strategies, reviewed, corrected and discussed the derivations and the statistical

832     aspects of the paper. M.A. Elzo designed the simulation, reviewed, corrected and discussed the

833     genetic aspects of the paper.

834     **Acknowledgments**

843 **References**

844 Bernardo, R., Yu, J. (2007). Prospects for Genomewide Selection for Quantitative Traits in Maize.

845     *Crop Science*, 47, 1082-1090.

846 Casella, G., George, E.I. (1992). Explaining the Gibbs Sampler. *The American Statistician*, 46(3),

847     167-174.

848 Casella, G., Berger, R. (2002). *Statistical Inference* (2nd ed.). Duxbury, Pacific Grove, CA, USA.

849 Chen, L., Li, C., Miller, S., Schenkel, F. (2014). Multi-population genomic prediction using a multi-

850     task Bayesian learning model. *BMC Genetics*, 15:53.

851 de los Campos, G., Gianola D., Allison, D.B. (2010). Predicting genetic predisposition in humans:

852     the promise of whole-genome markers. *Nature Reviews Genetics*, 11, 880-886.

853 de los Campos, G., Veturi, Y., Vázquez, A.I., Lehermeier, C., Pérez-Rodríguez, P. (2015a).

854     Incorporating genetic heterogeneity in whole-genome regressions using interactions. Journal

855     of Agricultural, Biological and Environmental Statistics, 20(4), 467-490.

856 de los Campos, G., Sorensen, D., Gianola, D. (2015b). Genomic Heritability: What Is It? *PLOS*

857     *Genetics*, 11(5), e1005048.

858 de Roos, A.P.W., Hayes, B.J., Goddard, M.E. (2009). Reliability of Genomic Predictions Across

859     Multiple Populations. *Genetics*, 183, 1545-1553.

860 Desta, Z.A., Ortiz, R. (2014). Genomic selection: genome-wide prediction in plant improvement.

861     *Trends in Plant Science*, 19(9), 592-601.

862 Falconer, D.S., Mackay, T.F.C. (1996). *Introduction to Quantitative Genetics* (4th ed.). Longmans
863     Green, Harlow, UK.

864 Gelman, A., Carlin, J.B., Stern, H., Dunson, D.B., Vehtari, A., Rubin, D.B. (2013). *Bayesian Data*
865     *Analysis* (3rd ed.). Chapman and Hall/CRC, Boca Raton, FL, USA.

866 Gianola, D., Perez-Encizo, M., Toro, M.A. (2003). On Marker-Assisted Prediction of Genetic Value:
867     Beyond the Ridge. *Genetics*, 163, 347-365.

868 Gianola, D., de los Campos, G., Hill, W.G., Manfredi, E., Fernando, R.L. (2009) Additive genetic
869     variability and the Bayesian alphabet. *Genetics,* 183, 347-363.

870 Gianola, D., Simianer, H., Qanbari, S. (2010) A two-step method for detecting selection signatures
871     using genetic markers. *Genetic Research Cambridge*, 92(2), 141-155.

872 Gianola, D. (2013). Priors in whole-genome regression: The Bayesian alphabet returns. *Genetics,*
873     194, 573-596.

874 Gianola, D., de los Campos, G., Toro, M.A., Naya, H., Schön, C.C., Sorensen, D. (2015). Do
875     Molecular Markers Inform About Pleiotropy? *Genetics*, 201, 23-29.

876 Gianola, D., Rosa, G. (2015). One Hundred Years of Statistical Developments in Animal Breeding.
877     *Annual Reviews of Animal Biosciences*, 3, 19-56.

878 Goddard, M.E., Hayes, B.J. (2007). Genomic Selection. *Journal of Animal Breeding and Genetics*,
879     124, 323-330.

880 Goddard, M.E. (2009). Genomic selection: prediction of accuracy and maximization of long term
881     response. *Genetica*, 136, 245-257.

882 Guttmacher, A.E., Collins, F.S. (2002).Genomic medicine- a primer. *The New England Journal of*
883     *Medicine*, 347, 1512-1520.

884 Hayes, B.J., Bowman, P.J., Chamberlain, A.J., Goddard, M.E. (2009). Invited review: Genomic
885     selection in dairy cattle: Progress and challenges. *Journal of Dairy Science*, 92, 433-443.

886  Hill, W.G. (1984). On selection among groups with heterogeneous variance. *Animal Production*,
887      39(3), 473-477.

888  Hill, W.H. (2014). Applications of Population Genetics to Animal Breeding, from Wright, Fisher and
889      Lush to Genomic Prediction. *Genetics*, 196, 1-16.

890  Henderson, C.R. (1974). Use of all relatives in intraherd prediction of breeding values and producing
891      abilities. *Journal of Daity science, 58(12),* 1910-1916.

892  Huang, H., Windig, J.J., Vereijken, A., Calus, M.P.L. (2014). Genomic prediction based on data from
893      three layer lines using non-linear regression models. *Genetics Selection Evolution*, 46:75.

894  Jia, Y., Jannink, J.L. (2012) Multiple-Trait Genomic Selection Methods Increase Genetic Value
895      Prediction Accuracy. *Genetics*, 192, 1513-1522.

896  Karoui, S., Carabaño, M.J., Díaz, C., Legarra, A. (2012). Joint genomic evaluation of French dairy
897      cattle breeds using multiple-trait models. *Genetics Selection Evolution*, 44:39.

898  Kennedy, B.W., Schaeffer, L.R., Sorensen, D.A. (1988). Genetic properties of animal models.
899      *Journal of Dairy Science, 71(2),* 17-26.

900  Lehermeir, C., Schon, C., de los Campos, G. (2015). Assessment of genetic heterogeneity in
901      structured plant populations using multivariate whole-genome regression models. *Genetics*,
902      201, 323-337.

903  Li, Y., Willer, C., Sanna, S., Abecasis, G. (2009). Genotype imputation. *Annual Review of Genomics*
904      *and Human Genetics*, *10*, 387-406.

905  Lund, M.S., de Roos, A.P.W., de Vries, A.G., Druet, T., Ducrocq, V., Fritz, S., Guillaume, F.,
906      Guldbrandtsen, B., Liu, Z., Reents, R., Schrooten, C., Seefried, F., Su, J. (2011). A common
907      reference population from four European Holstein populations increases reliability of
908      genomic predictions. *Genetics Selection Evolution*, 43:43.

909  Lynch, M., Walsh, E. (1998). *Genetics and Analysis of Quantitative Traits*. Sinauer associates Inc.,
910      Sunderland, MA, USA.

911   Maier, R., Moser, G., Chen, G.B., Ripke, S., Cross-Disorder Working Group of the Psychiatric

912        Genomics Consortium, Coryell, W., Potash, J.B., Scheftner, W.A., Shi, J., Weissman, M.M.,

913        Hultman, C.M., Lande´n, M., Levinson, D.F., Kendler, K.S., Smoller, J.W., Wray, N.R., Lee,

914        S.H. (2015). Joint Analysis of Psychiatric Disorders Increases Accuracy of Risk Prediction

915        for Schizophrenia, Bipolar Disorder, and Major Depressive Disorder. *The American Journal*

916        *of Human Genetics*, 96, 283-294.

917   Makgahlela, M.L., Mantysaari, E.A., Stranden, I., Koivula, M., Nielsen, U.S., Sillanpaa, M.J., Juga,

918        J. (2013). Acroos breed multi-trait random regression genomic predictions in the Nordic Red

919        dairy cattle. *Journal of Animal Breeding and Genetics*, 130, 10-19.

920   Martínez, C.A., Khare, K., Elzo, M.A. (2015). On the Bayesness, minimaxity and admissibility of

921        point estimators of allelic frequencies. *Journal of Theoretical Biology*, 383, 106-115.

922   Meuwissen, T.H.E., Hayes B.J., Goddard, M.E. (2001) Prediction of total genetic value using

923        genome-wide dense marker maps. *Genetics,* 157,1819-1829.

924   Olson, K.M., VanRaden, P.M., Tooker, M.E. (2012). Multibreed genomic evaluations using purebred

925        Holsteins, Jerseys, and Brown Swiss. *Journal of Dairy Science*, 95, 5378-5383.

926   R Core Team (2015). R: A language and environment for statistical computing. R foundation for

927        statistical computing, Vienna, Austria. *URL https://www.R-project.org/.*

928   Robert, C.P., Casella, G. (2010). *Introducing Monte Carlo Methods with R*. Springer, New York,

929        NY, USA.

930   Sargolzaei, M., Schenkel, F.S. (2013). *QMSim User's Guide Version 1.10*. Centre for Genetic

931        Improvement of Livestock, Department of Animal and Poultry Science,    University    of

932        Guelph, Guelph, Canada.

933   Spiegelhalter, D.J., Best, N.G., Carlin, B.P., van der Linde, A. (2002). Bayesian measures of model

934        complexity and fit (with discussion). *Journal of the Royal Statistical Society Series B*, 64,

935        583–639.

936    Spiegelhalter, D.J., Best, N.G., Carlin, B.P., van der Linde, A. (2014). The deviance information

937         criterion: 12 years on. *Journal of the Royal Statistical Society Series B*, 76, 485–493.

938    van den Berg, S., Calus, M.P.L., Meuwissen, T.H.E., Wientjes, Y.C.J.  (2015). Across population

939         genomic prediction scenarios in which Bayesian variable selection outperforms GBLUP.

940         *BMC Genetics*, 16:416.

941    Wientjes, Y.C.J., Veerkamp, R.F., Bijma, P., Bovenhuis, H., Schrooten, C., Calus, M.P.L. (2015)

942         Empirical and deterministic accuracies of across-population genomic prediction. *Genetics*

943         *Selection Evolution*, 47:5.

944    Wright, S.  (1930). Evolution in Mendelian populations. *Genetics,* 16, 98-159.

945    Wright, S.  (1937). The distribution of genetic frequencies in populations. *Genetics,* 23, 307-320.

946    Yang, W., Tempelman, R.J. (2012). A Bayesian Antedependence Model for Whole Genome

947         Prediction. *Genetics*, 190, 1491-1501.

948    Zhong, S.,  Dekkers, J.C.M., Fernando, R.L., Jannink, J.K. (2009). Factors Affecting Accuracy From

949         Genomic Selection in Populations Derived From Multiple Inbred Lines: A Barley Case

950         Study. *Genetics*, 182, 355-364.

951
952
953
954
955
956
957
958
959
960
961
962
963
964
965

## Appendix A: Conditional pmf of genotypes given parental genotypes, joint posteriors, full conditionals and details of some derivations

**Table A.1** Conditional pmf of genotypes at locus $j$ given the parental genotypes

| Parental genotypes | | Corresponding random variables | | $\pi\left(w_{ij}\|w_{S_ij}, w_{D_ij}\right) = \Pr(w_{ij} = x\|w_{S_ij} = k, w_{D_ij} = z)$ $x, k, z \in \{-1,0,1\}$ | | |
|---|---|---|---|---|---|---|
| Parent 1 | Parent 2 | $w_{S_ij}$ | $w_{D_ij}$ | $\pi\left(-1\|w_{S_ij}, w_{D_ij}\right)$ | $\pi\left(0\|w_{S_ij}, w_{D_ij}\right)$ | $\pi\left(1\|w_{S_ij}, w_{D_ij}\right)$ |
| AA | AA | -1 | -1 | 1 | 0 | 0 |
| AA(BB) | BB(AA) | -1(1) | 1(-1) | 0 | 1 | 0 |
| AB | AB | 0 | 0 | ¼ | ½ | ¼ |
| AA(AB) | AB(AA) | -1(0) | 0(-1) | ½ | ½ | 0 |
| BB(AB) | AB(BB) | 1(0) | 0(1) | 0 | ½ | ½ |
| BB | BB | 1 | 1 | 0 | 0 | 1 |

**Joint posteriors for Homogeneous marker effect covariance matrix model with homoscedastic residuals and Gaussian prior for $g$**

Weather $P^*$ is considered a parameter (some founders are genotyped) or a hyperparameter (none of the founders is genotyped) is not relevant when computing the joint posterior because in both cases its pdf is the same, thus it enters in the expression in the same way. Henceforth, it is assumed that vector $g$ and columns of matrix $W$ are ordered by marker unless otherwise indicated. Thus:

$$\pi(g, \sigma^2, W^N, G^0, P^*|y, W^o) \propto f(y|g, \sigma^2, W)\pi(g|G^0)\pi(G^0)\pi(\sigma^2)\pi(W|P^*)\pi(P^*)$$

$$\propto (\sigma^2)^{-\frac{n}{2}} \exp\left(\frac{-1}{2\sigma^2}(y - Wg)'(y - Wg)\right)$$

$$\times |G^0|^{-\frac{m}{2}} \exp\left(\frac{-1}{2}g'(I_m \otimes (G^0)^{-1})g\right)$$

$$\times |G^0|^{-\frac{1}{2}(a+S+1)} \exp\left(\frac{-1}{2}tr(\Sigma(G^0)^{-1})\right)$$

$$\times (\sigma^2)^{-\left(\frac{v}{2}+1\right)} \exp\left(\frac{-\tau^2}{2\sigma^2}\right)$$

$$\times \pi(W|P^*)\pi(P^*)$$

Where $\otimes$ represents the Kronecker product and $\pi(W|P^*)\pi(P^*) = \pi(W, P|r)$, when $r$ is assumed to be known and has the following form (see appendix A for details):

$$\pi(W, P|r) \propto$$

$$2^{n_H} \prod_{j=1}^{m} p_{(S+1)j}^{\alpha_{S+1}-1} \prod_{l=1}^{S} \left\{ \frac{1}{r_l^{2f_l}} p_{lj}^{n_l^{B_j}+\alpha_l-1}(r_l - p_{lj})^{n_l^{A_j}} \prod_{i'=f_l+1}^{n_l} \pi\left(w_{i'j}^l|w_{S_{i'j}}, w_{D_{i'j}}\right) \right\}.$$

982  When $r$ is unknown, the only change is that expression $(p_{lj})^{n_l^{B_j}+\alpha_l-1}(r_l-p_{lj})^{n_l^{A_j}}$ has to be replaced

983  by $(p_{lj})^{n_l^{B_j}+\alpha_{lp}-1}(r_{lj}-p_{lj})^{n_l^{A_j}+\alpha_{lq}-1}$ and instead of $\pi(W,P|r), \pi(W\,|\,P^*)\pi(P^*)$ corresponds to

984  $\pi(W,P,Q), Q := (\boldsymbol{q}_1,\dots,\boldsymbol{q}_j)$.

985

986  **Joint posteriors for Heterogeneous marker effect covariance matrix model with homoscedastic**
987  **residuals and Gaussian prior for $\boldsymbol{g}$**

988

$$\pi(\boldsymbol{g},\sigma^2,W^N,G,P|\boldsymbol{y},W^\sigma) \propto f(\boldsymbol{y}|\boldsymbol{g},\sigma^2,W)\pi(\boldsymbol{g}|G)\pi(G)\pi(\sigma^2)\pi(W|P)\pi(P)$$

$$\propto (\sigma^2)^{-\frac{n}{2}}\exp\left(\frac{-1}{2\sigma^2}(\boldsymbol{y}-W\boldsymbol{g})'(\boldsymbol{y}-W\boldsymbol{g})\right)$$

$$\times |G^0|^{-\frac{m}{2}}\exp\left(\frac{-1}{2}\boldsymbol{g}'(I_m\otimes(G^0)^{-1})\boldsymbol{g}\right)$$

$$\times \prod_{j=1}^{m}\left\{|G_j|^{-\frac{1}{2}(a+\mathcal{S}+1)}\right\}\exp\left(\frac{-1}{2}\sum_{j=1}^{m}tr(\Sigma G_j^{-1})\right)$$

$$\times (\sigma^2)^{-\left(\frac{v}{2}+1\right)}\exp\left(\frac{-\tau^2}{2\sigma^2}\right)$$

$$\times \pi(W|P^*)\pi(P^*).$$

989

990

991  **Marginal prior distribution of marker effects**

992

993  *Homogeneous marker effect covariance matrix models*

994

$$\pi(\boldsymbol{g}) \propto \int_{\mathcal{P}_\mathcal{S}^+} \pi(\boldsymbol{g}|G^0)\pi(G^0)dG^0$$

$$\propto \int_{\mathcal{P}_\mathcal{S}^+} \exp\left(\frac{-1}{2}tr\left(\left(\Sigma+\sum_{j=1}^{m}\boldsymbol{g}_j\boldsymbol{g}_j'\right)(G^0)^{-1}\right)\right)|G^0|^{-\frac{1}{2}(a+\mathcal{S}+m+1)}dG^0$$

995

996  the expression $tr\left(\left(\Sigma+\sum_{j=1}^{m}\boldsymbol{g}_j\boldsymbol{g}_j'\right)(G^0)^{-1}\right)$ comes from adding terms $\boldsymbol{g}'(I_m\otimes(G^0)^{-1})\boldsymbol{g}$ coming
997  from $\pi(\boldsymbol{g}|G^0)$ and $tr(\Sigma(G^0)^{-1})$ coming from $\pi(G^0)$. The equality is shown using properties of the
998  $tr(\cdot)$ operator as follows:

$$\boldsymbol{g}'(I_m\otimes(G^0)^{-1})\boldsymbol{g} = tr(\boldsymbol{g}'(I_m\otimes(G^0)^{-1})\boldsymbol{g})$$

$$= tr\left((\boldsymbol{g}_1 \cdots \boldsymbol{g}_m)\begin{pmatrix}(G^0)^{-1} & & \\ & \ddots & \\ & & (G^0)^{-1}\end{pmatrix}\begin{pmatrix}\boldsymbol{g}_1 \\ \vdots \\ \boldsymbol{g}_m\end{pmatrix}\right)$$

$$= tr\left(\sum_{j=1}^{m} \boldsymbol{g}_j'(G^0)^{-1}\boldsymbol{g}_j\right)$$

$$= tr\left(\sum_{j=1}^{m} \boldsymbol{g}_j\boldsymbol{g}_j'(G^0)^{-1}\right),$$

moreover, since $tr(\boldsymbol{g}'(I_m\otimes(G^0)^{-1})\boldsymbol{g}) = tr(\boldsymbol{g}\boldsymbol{g}'(I_m\otimes(G^0)^{-1}))$, it follows that:

$$\boldsymbol{g}'(I_m\otimes(G^0)^{-1})\boldsymbol{g} + tr(\boldsymbol{\Sigma}(G^0)^{-1}) = tr\left(\boldsymbol{\Sigma}(G^0)^{-1} + \boldsymbol{g}\boldsymbol{g}'(I_m\otimes(G^0)^{-1})\right)$$

$$= tr\left(\boldsymbol{\Sigma}(G^0)^{-1} + \sum_{j=1}^{m}\boldsymbol{g}_j\boldsymbol{g}_j'(G^0)^{-1}\right)$$

$$= tr\left(\left(\boldsymbol{\Sigma} + \sum_{j=1}^{m}\boldsymbol{g}_j\boldsymbol{g}_j'\right)(G^0)^{-1}\right).$$

Using this, it follows that:

$$\pi(\boldsymbol{g}) \propto \int_{\mathcal{P}_\mathcal{S}^+} \exp\left(\frac{-1}{2}tr\left(\left(\boldsymbol{\Sigma} + \sum_{j=1}^{m}\boldsymbol{g}_j\boldsymbol{g}_j'\right)(G^0)^{-1}\right)\right)|G^0|^{-\frac{1}{2}(a+\mathcal{S}+m+1)}dG^0$$

$$= \frac{2^{(a+m)\mathcal{S}/2}\Gamma_\mathcal{S}\left(\frac{a+m}{2}\right)}{\left|\boldsymbol{\Sigma} + \sum_{j=1}^{m}\boldsymbol{g}_j\boldsymbol{g}_j'\right|^{\left(\frac{a+m}{2}\right)}}$$

This result easily follows because we are integrating the kernel of an inverse Wishart density with parameters $\left(\boldsymbol{\Sigma} + \sum_{j=1}^{m}\boldsymbol{g}_j\boldsymbol{g}_j', a + m\right)$.

*Heterogeneous marker effect covariance matrix model*

For this model:

$$\pi(\boldsymbol{g}) \propto \prod_{j=1}^{m} \int_{\mathcal{P}_\mathcal{S}^+} |G_j|^{-\frac{1}{2}(a+\mathcal{S}+2)} \exp\left(\frac{-1}{2}tr\left(\left(\boldsymbol{\Sigma} + \boldsymbol{g}_j\boldsymbol{g}_j'\right)(G_j)^{-1}\right)\right)dG_j$$

$$= \prod_{j=1}^{m} \frac{2^{(a+1)\mathcal{S}/2}\Gamma_\mathcal{S}\left(\frac{a+1}{2}\right)}{\left|\boldsymbol{\Sigma} + \boldsymbol{g}_j\boldsymbol{g}_j'\right|^{\left(\frac{a+1}{2}\right)}}$$

$$= \frac{2^{(a+1)m\mathcal{S}/2}\left(\Gamma_{\mathcal{S}}\left(\frac{a+1}{2}\right)\right)^m}{\prod_{j=1}^{m}|\boldsymbol{\Sigma} + \boldsymbol{g}_j\boldsymbol{g}_j'|^{\left(\frac{a+1}{2}\right)}},$$

using the results for determinants of partitioned matrices this expression can be written as:

$$\frac{2^{(a+1)m\mathcal{S}/2}\left(\Gamma_{\mathcal{S}}\left(\frac{a+1}{2}\right)\right)^m}{\prod_{j=1}^{m}\left(|\boldsymbol{\Sigma}|\, 1 + \boldsymbol{g}_j'\boldsymbol{\Sigma}^{-1}\boldsymbol{g}_j\right)^{\left(\frac{a+1}{2}\right)}} \propto \frac{1}{\prod_{j=1}^{m}\left(1 + \frac{1}{a+1-\mathcal{S}}\boldsymbol{g}_j'\boldsymbol{\Sigma}_*^{-1}\boldsymbol{g}_j\right)^{\left(\frac{a+1}{2}\right)}},$$

where $\boldsymbol{\Sigma}_* = \frac{1}{a+1-\mathcal{S}}\boldsymbol{\Sigma}$. This is the product of multivariate t distributions with scale matrix $\boldsymbol{\Sigma}_*$ and

degrees of freedom $a + 1 - \mathcal{S}$.

1014

**Details on the form of $\pi(W, P^*)$, $r$ known**

1016

$$\pi(W|P^*)\pi(P^*) = \pi(W|P,\boldsymbol{r})\pi(P|\boldsymbol{r})$$

$$= \pi(W|P,\boldsymbol{r})\prod_{j=1}^{m}\pi(p_j|\boldsymbol{r})$$

$$= \frac{2^{n^H}}{c}\prod_{j=1}^{m}\prod_{l=1}^{\mathcal{S}}\left\{\frac{1}{r_l^{2f_l}}p_{lj}^{n_l^{Bj}}(r_l - p_{lj})^{n_l^{Aj}}\prod_{i'=f_l+1}^{n_l}\pi\left(w_{i'j}^l|w_{S_{i'j}}, w_{D_{i'j}}\right)\right\}\prod_{j=1}^{m}\pi(p_j|\boldsymbol{r})$$

$$\propto \frac{2^{n^H}}{c}\prod_{j=1}^{m}\prod_{l=1}^{\mathcal{S}}\left\{\frac{1}{r_l^{2f_l}}p_{lj}^{n_l^{Bj}}(r_l - p_{lj})^{n_l^{Aj}}\prod_{i'=f_l+1}^{n_l}\pi\left(w_{i'j}^l|w_{S_{i'j}}, w_{D_{i'j}}\right)\right\} \times \prod_{j=1}^{m}p_{(\mathcal{S}+1)j}^{\alpha_{\mathcal{S}+1}-1}\prod_{l=1}^{\mathcal{S}}\left(\frac{p_{lj}}{r_l}\right)^{\alpha_l-1}$$

$$\propto \frac{2^{n^H}}{c}\prod_{j=1}^{m}p_{(\mathcal{S}+1)j}^{\alpha_{\mathcal{S}+1}-1}\prod_{l=1}^{\mathcal{S}}\left\{\frac{1}{r_l^{2f_l}}p_{lj}^{n_l^{Bj}+\alpha_l-1}(r_l - p_{lj})^{n_l^{Aj}}\prod_{i'=f_l+1}^{n_l}\pi\left(w_{i'j}^l|w_{S_{i'}}, w_{D_{i'}}\right)\right\}$$

1017

$p_{(\mathcal{S}+1)j} = 1 - \sum_{l=1}^{\mathcal{S}}p_{lj}$, for each $j$, $\boldsymbol{g}_j \in \mathbb{R}^{\mathcal{S}}$ corresponds to the subvector of $\boldsymbol{g}$ containing the effects

of marker $j$ in each one of the $\mathcal{S}$ subpopulations and $\otimes$ represents the Kronecker product. Analogous

steps lead to the form of $\pi(W, P^*)$ when $\boldsymbol{r}$ is unknown.

1021

**Full conditionals**

1023

*Full conditionals for models with heteroscedastic residuals*

In this case:

$$f(\boldsymbol{y}|W, \boldsymbol{g}, R) \propto |V|^{-1/2}\exp\left(-\frac{1}{2}(\boldsymbol{y} - W\boldsymbol{g})'V^{-1}(\boldsymbol{y} - W\boldsymbol{g})\right)$$

$$= \prod_{l=1}^{\mathcal{S}}(\sigma_l^2)^{-n_l/2}\exp\left(-\frac{1}{2\sigma_l^2}(\boldsymbol{y}_l - W_l\boldsymbol{g}_l)'(\boldsymbol{y}_l - W_l\boldsymbol{g}_l)\right).$$

In addition

$$\pi(R) \propto \prod_{l=1}^{\mathcal{S}} (\sigma_l^2)^{-(v/2+1)} \exp\left(-\frac{\tau^2}{2\sigma_l^2}\right).$$

1027    In the following, only the full conditionals that change with respect to the homoscedastic models are
1028    presented. For the homogeneous marker effect covariance matrix model with multivariate normal
1029    prior the full conditionals that change are:

$$\pi(\boldsymbol{g}|Else) = MVN((W'V^{-1}W + G^{-1})^{-1}W'V^{-1}\boldsymbol{y}, (W'V^{-1}W + G^{-1})^{-1})$$

1030    where $G^{-1} = (G^0)^{-1} \otimes I$.

$$\pi(R|Else) = \prod_{l=1}^{\mathcal{S}} IG\left(\frac{v + n_l}{2}, \frac{\tau^2 + (\boldsymbol{y}_l - W_l\boldsymbol{g}_l)'(\boldsymbol{y}_l - W_l\boldsymbol{g}_l)}{2}\right).$$

1031    To define $\pi(W^N|Else)$ the partitions defined in section 2.2.1 are done for each subpopulation.

$$\pi(W^N|Else) \propto \pi^+(W|P^*) \prod_{l=1}^{\mathcal{S}} \exp\left(\frac{-1}{2\sigma_l^2}(-2\boldsymbol{g}_l'W_l^{N\prime}\boldsymbol{y}_l^N + \boldsymbol{g}_l'W_l^{N\prime}W_l^N\boldsymbol{g}_l)\right)$$

$$\times \prod_{l=1}^{\mathcal{S}} \prod_{k=1}^{K} \exp\left(\frac{-1}{2\sigma_l^2} h\left(W_l^{M_k}, \boldsymbol{g}_l^{M_k}, \boldsymbol{y}_l^{M_k}\right)\right)$$

1032    where

$$h\left(W_l^{M_k}, \boldsymbol{g}_l^{M_k}, \boldsymbol{y}_l^{M_k}\right)$$
$$= 2\left(\boldsymbol{g}_l^{M_kN\prime}W_l^{M_kN\prime}W_l^{M_k\sigma}\boldsymbol{g}_l^{M_k\sigma} - \boldsymbol{g}_l^{M_kN\prime\prime}W_l^{M_kN\prime}\boldsymbol{y}_l^{M_k}\right) + \boldsymbol{g}_l^{M_kN\prime}W_l^{M_kN\prime}W_l^{M_kN}\boldsymbol{g}_l^{M_kN}.$$

1033

1034    For the heterogeneous marker effect covariance matrix model with multivariate Gaussian prior for $\boldsymbol{g}$:
1035

$$\boldsymbol{g}|Else \sim MVN((W'V^{-1}W + G^{-1})^{-1}W'V^{-1}\boldsymbol{y}, (W'V^{-1}W + G^{-1})^{-1})$$

1036

1037    where $G^{-1} = Block\ Diag\ (G_1^{-1}, \ldots, G_{\mathcal{S}}^{-1})$.

1038
1039                  **Appendix B: Details on data simulation**
1040

1041    For phenotype one (dataset 1), in a first stage three preliminary subpopulations were simulated by
1042    selecting individuals from the historical population. Numbers of individuals and criteria to select
1043    them were the following. In preliminary subpopulation 1, ten males and 250 females with the lowest
1044    true breeding values, in preliminary subpopulation 2, five males and 200 females with the highest
1045    phenotypes and in preliminary subpopulation 3, 50 males and 500 females randomly chosen. Then,
1046    selection criteria and mating design to create new generations were: lowest phenotypes and positive
1047    assortative in preliminary subpopulation 1, highest phenotypic values and random for preliminary
1048    subpopulation 2, and random and random for preliminary subpopulation 3. Positive assortative
1049    means that individuals are mated looking for similarity, while negative assortative means that
1050    individuals are mated looking for dissimilarity, where (di)similarity can be defined in terms of
1051    phenotypes, true or predicted breeding values (Sargolzaei and Schenkel, 2013). The numbers of
1052    simulated generations were four, two, and three respectively. Subsequently, two more subpopulations
1053    hereinafter referred to as subpopulations one and two were simulated as follows. Eighteen males and

1054 100 females from the fourth generation of the first subpopulation, two males and 40 females from the
1055 second generation of the second subpopulation, and eight males and 40 females from the third
1056 generation of the third subpopulation were chosen to create the subpopulation one. Ten females from
1057 generation three of preliminary subpopulation one, 20 males and 60 females from generation two of
1058 preliminary subpopulation two, and 20 females from generation two of preliminary subpopulation
1059 three were chosen to generate subpopulation two. Generations zero and one of preliminary
1060 subpopulation three were used to define subpopulation three. For the second phenotype (dataset 2)
1061 the two subpopulations were simulated by choosing individuals from the historical subpopulation
1062 based on different criteria and mating them according to different systems and selection criteria for
1063 two generations.
1064 In each case, a single pedigree was simulated which allowed individuals from a given subpopulation
1065 to be parents of individuals from another subpopulation. This mimics what happens in certain
1066 populations like animal populations when using semen or oocytes from individuals from a different
1067 subpopulation (e.g., country) to produce a new generation of a given subpopulation. The number of
1068 alleles per QTL was two, three and four; these numbers were randomly assigned using a uniform
1069 distribution. QTL were evenly allocated across the genome as well as SNP markers.
1070 In both datasets, additive QTL effects were scaled such that QTL effects and heritabilities were
1071 different in each subpopulation. Within a given subpopulation, all QTL allelic effects were scaled by
1072 the same factor. Markers with minor allele frequencies smaller than 0.05 were excluded from the
1073 analysis.
1074
1075 **Appendix C: Conditional pmf of genotypes at locus *j* given one parental genotype and allelic**
1076 **frequencies**
1077
1078 The following table shows $\pi\left(w_{ij}\middle|w_{Pa_{ij}},\boldsymbol{p}_j^*\right) = \Pr(w_{ij} = x|w_{Pa_{ij}} = z,\boldsymbol{p}_j^*), x, z \in \{-1,0,1\}$, where
1079 $w_{Pa_{ij}}$ is the genotype of the known parent of individual $i$ for marker locus $j$. If the subpopulation to
1080 which the unknown parent pertains is known to be subpopulation $l$ then $\pi\left(w_{ij}\middle|w_{Pa_{ij}},\boldsymbol{p}_j^*\right)$ has the
1081 following form:
1082

| Known parental genotype | $w_{Pa_{ij}}$ | $\pi\left(-1\middle|w_{Pa_{ij}},\boldsymbol{p}_j^*\right)$ | $\pi\left(0\middle|w_{Pa_{ij}},\boldsymbol{p}_j^*\right)$ | $\pi\left(1\middle|w_{Pa_{ij}},\boldsymbol{p}_j^*\right)$ |
|---|---|---|---|---|
| AA | -1 | $1 - p_{lj}^*$ | $p_{lj}^*$ | 0 |
| AB | 0 | $\left(1 - p_{lj}^*\right)/2$ | $1/2$ | $p_{lj}^*/2$ |
| BB | 1 | 0 | $1 - p_{lj}^*$ | $p_{lj}^*$ |

1083
1084 If no information about the unknown parent is available, one pragmatic solution is to assume that the
1085 probabilities of inherit a given allele are dictated by the unweighted average of allelic frequencies
1086 across subpopulations (for the full models). If $\bar{p}_j^*$ represents that average reference allele frequency
1087 for marker locus $j$ then the conditional probabilities are same as in the previous table with $p_{lj}^*$
1088 replaced by $\bar{p}_j^*$. Of course, the lack of knowledge of the origin of the unknown parent is not an issue
1089 for null models.

# Appendix D: Toy example of the joint pmf of two half sib founders

In this case the common parent is the relevant common ancestor. This individual is denoted with
number 3. Suppose that individuals 1, 2 and 3 belong to population $l$. For simplicity we focus on a
single marker, thus the subindex associated with marker is ignored. Then:

$$\pi(w_1, w_2 | P^*) = \sum_{k \in \{-1,0,1\}} \pi(w_1 | w_3 = k, p_l^*) \pi(w_2 | w_3 = k, p_l^*) \pi(w_3 = k | p_l^*).$$

This summation is done for every one of the 9 combinations of genotypes of individuals 1 and 2. The
following table displays the conditional probabilities $(w_1, w_2 | w_3 = k, p_l^*)$.

| Genotype of 3 | Genotype of 2 | Genotype of 1 | | |
|---|---|---|---|---|
| | | AA | AB | BB |
| AA | AA | $(1 - p_l^*)^2$ | $p_l^*(1 - p_l^*)$ | 0 |
| | AB | $p_l^*(1 - p_l^*)$ | $p_l^{*2}$ | 0 |
| | BB | 0 | 0 | 0 |
| AB | AA | $(1 - p_l^*)^2/4$ | $(1 - p_l^*)/4$ | $p_l^*(1 - p_l^*)/4$ |
| | AB | $(1 - p_l^*)/4$ | $1/4$ | $p_l^*/4$ |
| | BB | $p_l^*(1 - p_l^*)/4$ | $p_l^*/4$ | $p_l^{*2}/4$ |
| BB | AA | 0 | 0 | 0 |
| | AB | 0 | $(1 - p_l^*)^2$ | $p_l^*(1 - p_l^*)$ |
| | BB | 0 | $p_l^*(1 - p_l^*)$ | $p_l^{*2}$ |

The following table presents the joint pmf of individuals 1 and 2 conditional on allelic frequencies

| Genotype of 2 | Genotype of 1 | | |
|---|---|---|---|
| | AA | AB | BB |
| AA | $(1 - p_l^*)^3 \left(1 - \frac{p_l^*}{2}\right)$ | $(1 - p_l^*)^2 p_l^* \left(\frac{3}{2} - p_l^*\right)$ | $\frac{\left(p_l^*(1 - p_l^*)\right)^2}{2}$ |
| AB | $(1 - p_l^*)^2 p_l^* \left(\frac{3}{2} - p_l^*\right)$ | $p_l^*(1 - p_l^*)\left(2p_l^*(1 - p_l^*) + \frac{1}{2}\right)$ | $p_l^{*2}(1 - p_l^*)\left(p_l^* + \frac{1}{2}\right)$ |
| BB | $\frac{\left(p_l^*(1 - p_l^*)\right)^2}{2}$ | $p_l^{*2}(1 - p_l^*)\left(p_l^* + \frac{1}{2}\right)$ | $p_l^{*3}\left(\frac{p_l^* + 1}{2}\right)$ |