Available online at [www.sciencedirect.com](http://www.sciencedirect.com)

ScienceDirect



# Genomic selection methods for crop improvement: Current status and prospects

Xin Wang<sup>a,b,c</sup>, Yang Xu<sup>a</sup>, Zhongli Hu<sup>c</sup>, Chenwu Xu<sup>a,\*</sup>

<sup>a</sup>Jiangsu Provincial Key Laboratory of Crop Genetics and Physiology/Co-Innovation Center for Modern Production Technology of Grain Crops/Key Laboratory of Plant Functional Genomics of Ministry of Education, Yangzhou University, Yangzhou 225009, Jiangsu, China

<sup>b</sup>College of Information Engineering, Yangzhou University, Yangzhou 225009, Jiangsu, China

<sup>c</sup>State Key Laboratory of Hybrid Rice, College of Life Sciences, Wuhan University, Wuhan 430072, Hubei, China

## ARTICLE INFO

### Article history:

Received 4 February 2018

Received in revised form 25 March 2018

Accepted 9 April 2018

Available online 15 April 2018

### Keywords:

Genomic selection

Prediction

Accuracy

Crop

## ABSTRACT

With marker and phenotype information from observed populations, genomic selection (GS) can be used to establish associations between markers and phenotypes. It aims to use genome-wide markers to estimate the effects of all loci and thereby predict the genetic values of untested populations, so as to achieve more comprehensive and reliable selection and to accelerate genetic progress in crop breeding. GS models usually face the problem that the number of markers is much higher than the number of phenotypic observations. To overcome this issue and improve prediction accuracy, many models and algorithms, including GBLUP, Bayes, and machine learning have been employed for GS. As hot issues in GS research, the estimation of non-additive genetic effects and the combined analysis of multiple traits or multiple environments are also important for improving the accuracy of prediction. In recent years, crop breeding has taken advantage of the development of GS. The principles and characteristics of current popular GS methods and research progress in these methods for crop improvement are reviewed in this paper.

© 2018 Crop Science Society of China and Institute of Crop Science, CAAS. Production and hosting by Elsevier B.V. on behalf of KeAi Communications Co., Ltd. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

## Contents

|  |     |
|--|-----|
| 1. Introduction . . . . .  | 331 |
| 2. Outline of genomic selection in breeding . . . . .              | 331 |
| 3. Training population design for crop genomic selection . . . . . | 331 |
| 4. Prediction models in genomic selection . . . . .                | 332 |
| 4.1. Prediction methods for additive genetic effects . . . . .     | 332 |
| 4.1.1. RR-BLUP . . . . .   | 333 |
| 4.1.2. BayesA . . . . .  | 333 |
| 4.1.3. BayesB . . . . .  | 333 |
| 4.1.4. BayesC . . . . .  | 333 |

\* Corresponding author.

E-mail address: [cwxu@yzu.edu.cn](mailto:cwxu@yzu.edu.cn). (C. Xu).

Peer review under responsibility of Crop Science Society of China and Institute of Crop Science, CAAS.

|        |   |     |
|--------|---|-----|
| 4.1.5. | LASSO . . . . .   | 334 |
| 4.1.6. | GBLUP . . . . .   | 334 |
| 4.1.7. | Reproducing kernel Hilbert space . . . . .                    | 334 |
| 4.1.8. | Support vector machine . . . . .                              | 334 |
| 4.1.9. | Characteristics of general methods . . . . .                  | 335 |
| 5.     | Factors affecting the accuracy of genomic selection . . . . . | 335 |
| 6.     | Perspective . . . . .   | 337 |
|        | Acknowledgments . . . . .                                     | 338 |
|        | References . . . . .  | 338 |

## 1. Introduction

The key step in crop breeding is selection, and conventional breeding is based on phenotypic selection. Breeders choose good offspring using their experience and the observed phenotypes of crops, so as to achieve genetic improvement of target traits. Hazel and Lush [1] first proposed the selection index (SI) method, which uses a total score to select for multiple traits simultaneously. It is more efficient than selection for one trait at a time and can improve aggregate genetic gain. In 1970s, with the development of computer science, Henderson [2] proposed best linear unbiased prediction (BLUP), which became the most widely used genetic evaluation method. Since 1990s, advances in molecular genetic techniques have revealed widespread genetic variation in genomes. Large numbers of molecular markers are available, allowing breeders to use markers to assist in breeding. Marker-assisted selection (MAS) [3] has become a common means of molecular breeding, but it is suitable only for traits controlled by a small number of major genes. Most economic traits of crops are complex and affected by a large number of genes, each with small effect [4], and thus the application of MAS in breeding practice is not as successful as expected. Owing to this limitation of MAS, it is necessary to incorporate genome-wide markers. Genomic selection (GS) [5] is an upgraded form of MAS. It aims to use genome-wide markers to estimate the effects of all loci and thereby compute a genomic estimated breeding value (GEBV), so as to achieve more comprehensive and reliable selection. Because the objects of selection are no longer limited to the traits determined by a few major genes, GS opens up a promising research direction for molecular breeding and it has become a hot issue in recent quantitative genetics research. In this paper, the principles and characteristics of several popular GS methods and research progress in these methods for crop improvement are reviewed.

## 2. Outline of genomic selection in breeding

GS can establish associations between markers and phenotypes based on a training population. The quantitative trait locus (QTL) detection step is skipped in GS, and a prediction model of phenotypic traits and genome-wide markers is constructed. Using this model, the genetic effect values of unobserved individuals are predicted, avoiding the omission of some small-effect markers that would fail a significance test. Even if the effect of each marker is very small, a large amount of marker information covering the whole genome

still has the potential to explain all the genetic variance. To perform crop genomic prediction, a large number of loci should be genotyped. In recent years, DNA marker technology has developed rapidly. In addition to electrophoresis, molecular hybridization, and molecular markers based on PCR, there are molecular markers based on chip technology, such as Diversity Array Technology (DArT) and single-nucleotide polymorphisms (SNP). These technologies have laid a strong foundation for the application of GS in crop breeding.

Many researchers in GS, have proposed a series of models and algorithms. Early GS studies were performed mainly for animal genetic prediction. Especially in dairy cattle, accuracy has been increased substantially for young candidates from the application of GS [6,7]. As expected, the development of GS is also beneficial for plant breeding. The genetic gain of GS in maize is higher than that of MAS [8,9], and higher than that of the conventional pedigree breeding [10]. When applied to hybrid breeding of crops, GS is even more efficient because genotypes of hybrids can be inferred from their inbred parents, leading to lower cost in genotyping [11,12]. Xu et al. [13] used 278 randomly selected rice hybrids derived from 210 recombinant inbred lines (RIL) as a training set and predicted 21,945 potential hybrids. The average yield of the top 100 showed a 16% increase compared with the average yield of all potential hybrids. In wheat, Daetwyler et al. [14] applied genomic best linear unbiased prediction (GBLUP) and a Bayesian regression method to predict rust resistance in 206 wheat landraces from 32 countries. The landraces were genotyped for 5568 SNPs using an Illumina 9 K chip, and the prediction accuracy of three traits ranged from 0.27 to 0.44. In short, GS uses all markers as predictors to achieve assessment and selection in early generations, so as to reduce the time cost per cycle and shorten the generation interval. Furthermore, GS saves labor costs compared with conventional breeding. It may dramatically change the role of phenotyping, which in GS serves to update prediction models and not only to select lines. In this paper, we present several popular GS methods and future directions for GS research that could revolutionize crop breeding.

## 3. Training population design for crop genomic selection

Fig. 1 depicts a road map of GS for crops, in which population design plays a vital role. In typical crop breeding, various elite breeding lines from a germplasm pool can be sampled as the training set for GS. However, such a simple strategy may limit

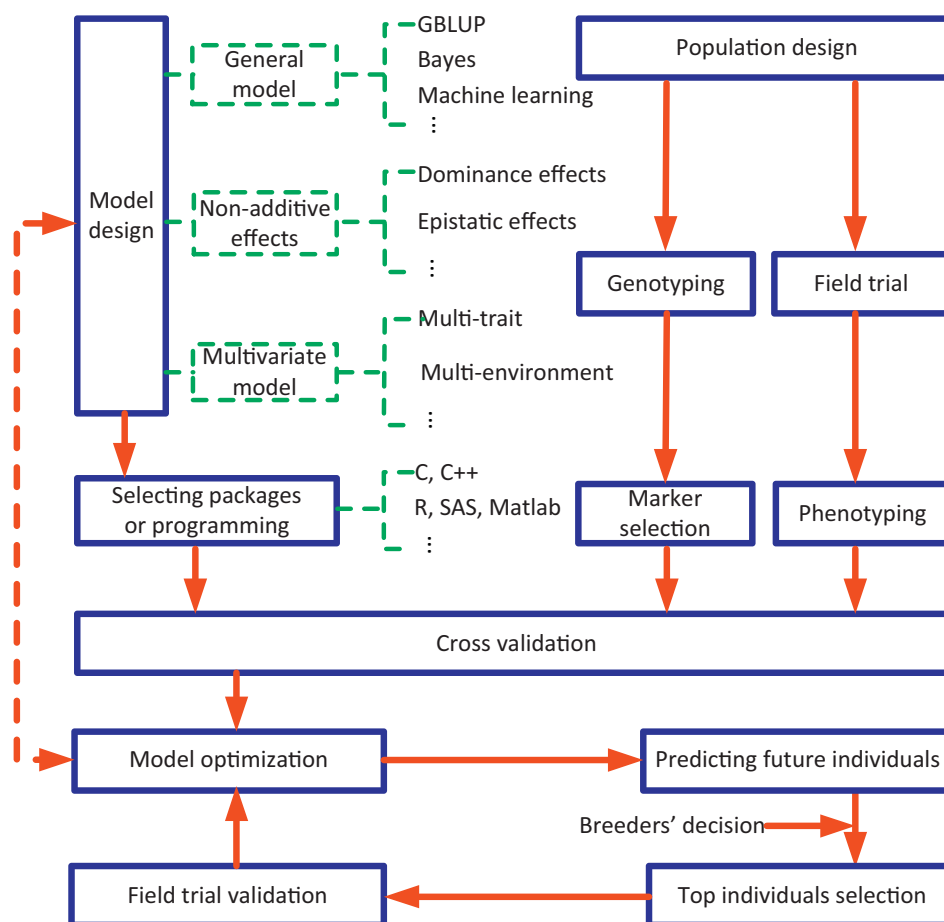


Fig. 1 – Road map of GS for crops.

the scope of prediction, and thus the benefit of GS cannot be significantly improved. It is known that hybrid breeding is successful at improving yield, especially for maize and rice. In this approach, breeders evaluate an inbred line not by its phenotype but by its potential to create superior hybrids, and thus selection of desirable hybrids is often a matter of trial and error in practice. GS is an important tool for facilitating hybrid breeding by obtaining better hybrids with fewer crosses.

In contrast to the half-sib families used in animal breeding, crop breeders often work with full-sib families. These are favorable for genomic prediction because biparental populations always have high linkage disequilibrium (LD) between markers and QTL. RIL, near-isogenic lines (NIL) and doubled haploid (DH) lines are major types of pure lines for crop breeding. In particular, DH technology has facilitated the generation of a large number of inbred lines which can not only accelerate the hybrid breeding process but also increase genetic gain.

Although biparental populations have been widely used in previous studies, their narrow genetic base may restrict the opportunity for achieving high success in breeding. Diallel cross designs can take advantage of a more comprehensive genetic resource, but are rarely used in the early stages of breeding programs owing to the large number of crosses. Alternatives are partial diallel and partial circular diallel designs, which are more feasible schemes than Diallel cross designs. Recently,

Wang et al. [15] predicted rice hybrid performance based on North Carolina mating design II (NC II). Using the rice data as a training set, performance of 15,115 potential crosses between five male sterile lines and the 3023 rice varieties in the 3000 rice genomes project [16] was predicted, showing the value of NC II scheme for GS and giving an example underlining the importance of population design.

## 4. Prediction models in genomic selection

In recent years, many methods have advanced GS, including general methods and their extensions (Fig. 2). General GS methods are based on additive models, and their accuracies may be different because they vary in their assumptions and algorithms with respect to the variances of complex traits. Incorporating non-additive effects or multiple variates, the general methods can be extended. The principles and characteristics of current popular GS methods are presented in this section.

### 4.1. Prediction methods for additive genetic effects

GS uses associations of a large number of markers across the whole genome with phenotypes to calculate accurate GEBVs of candidates for selection. However, in whole-genome

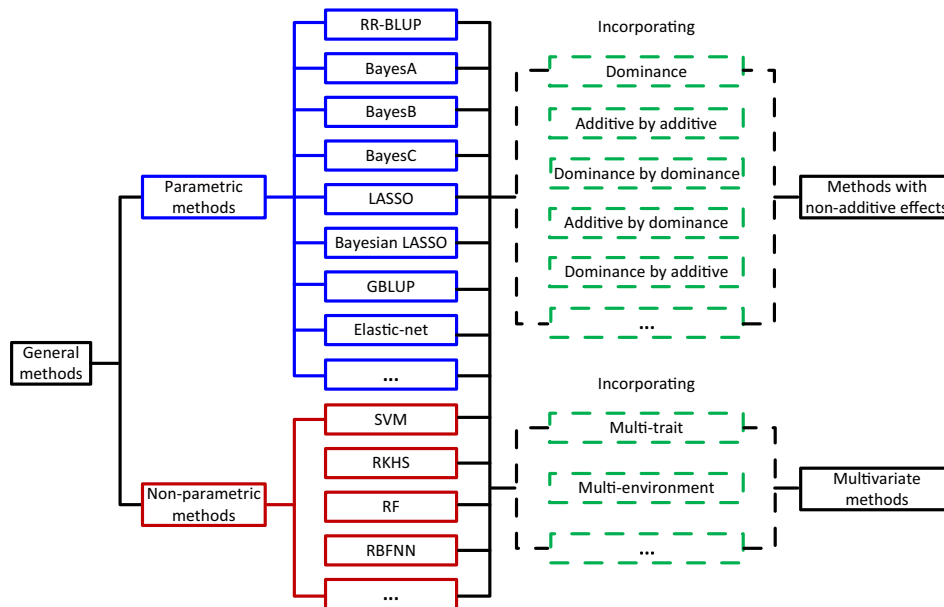


Fig. 2 – Classification of GS methods.

regression, the number of markers ( $k$ ) is usually much larger than the number of observations ( $n$ ). Degrees of freedom are insufficient to estimate all marker effects at the same time, a problem aggravated by multicollinearity [17]. If the effects of  $k$  markers are estimated simultaneously, the ordinary least squares method will be invalid. To address these issues, various methods such as GBLUP, Bayes and machine learning, have been employed for GS. The GBLUP and Bayesian methods treat the effects of markers as random effects, and the basic model is [18]:

$$y = X\beta + Z\alpha + \epsilon \quad (1)$$

where  $y$  is a vector of phenotypes;  $\beta$  is a vector of non-genetic fixed effects, for which a flat prior is often used;  $X$  is an incidence matrix for the fixed effects  $\beta$ ;  $\alpha$  is a vector of random regression coefficients of all the marker effects;  $Z$  is an  $n \times k$  genotypic matrix for markers; and  $\epsilon$  is a vector of residuals. Usually, a normal distribution with mean zero and covariance matrix  $R\sigma_\epsilon^2$  is used for  $\epsilon$ , where  $R$  is an identity matrix. Furthermore,  $\sigma_\epsilon^2$  has a scaled inverse chi-square distribution, that is,  $\sigma_\epsilon^2 \sim \chi^{-2}(\nu_\epsilon, S_\epsilon^2)$ . The alternative methods discussed here differ primarily in their specific prior used for  $\alpha$ .

#### 4.1.1. RR-BLUP

In ridge regression–best linear unbiased prediction (RR-BLUP),  $\alpha \sim N(0, I\sigma_\alpha^2)$  and  $\sigma_\alpha^2$  has a scaled inverse chi-square distribution; that is,  $\sigma_\alpha^2 \sim \chi^{-2}(\nu_\alpha, S_\alpha^2)$ . Then,  $\alpha$  has a marginal multivariate  $t$  distribution with mean zero, scale matrix  $S_\alpha^2 I$ , and degrees of freedom  $\nu_\alpha$ . The posterior of  $\sigma_\alpha^2$  is a scaled inverse chi-square distribution with scale  $\tilde{S}_\alpha^2 = (\sum_j \alpha_j^2 + \nu_\alpha S_\alpha^2) / \tilde{\nu}_\alpha$  ( $j = 1, \dots, k$ ) and degrees of freedom  $\tilde{\nu}_\alpha = \nu_\alpha + k$ . For RR-BLUP, a single effect variance is common to all loci, though this assumption of common variance does not imply that the effects of all markers are equal [9]. Xu et al. [19] pointed out that the

hypothesis is not consistent with fact in some cases and recommended an effect-specific prior variance.

#### 4.1.2. BayesA

Supposing a locus  $j$  is randomly sampled from among  $k$  loci, the marker effect at locus  $j$  has a univariate normal prior, with mean zero and locus-specific variance  $\sigma_j^2$ , which is assigned a scaled inverse chi-square prior in turn, that is  $\sigma_j^2 \sim \chi^{-2}(\nu_\alpha, S_\alpha^2)$ . It is important to note that each locus may have distinct variance, which is different from the case of RR-BLUP.  $\sigma_j^2$  has a posterior scaled inverse chi-square distribution with scale  $\tilde{S}_{\alpha_j}^2 = (\alpha_j^2 + \nu_\alpha S_\alpha^2) / \tilde{\nu}_\alpha$  and degrees of freedom  $\tilde{\nu}_\alpha = \nu_\alpha + 1$ . The unconditional distributions of the marker effects follow identical and independent univariate  $t$  distributions, each with mean zero.

#### 4.1.3. BayesB

BayesB employs a mixture distribution that includes a point of mass at zero and a univariate scaled  $t$  distribution. The  $t$  distribution is equivalent to a univariate normal with null mean and unknown locus-specific variance. For locus  $j$ , a Bernoulli variable  $\delta_j$  is allocated, which is 0 with probability  $\pi$  and 1 with probability  $(1 - \pi)$ . Thus, the effect of locus  $j$  can be calculated as  $\alpha_j = \xi_j \delta_j$ , where  $\xi_j \sim N(0, \sigma_j^2)$  and  $\sigma_j^2$  has a scaled inverse chi-square prior with scale parameter  $S_\alpha^2$  and degrees of freedom  $\nu_\alpha$ . The parameter  $\pi$  is always treated as a constant, which depends on the actual distribution of locus effects in real data analyses.

#### 4.1.4. BayesC

The assumption of BayesC is that each marker effect is zero with probability  $\pi$  and follows a univariate normal distribution with probability  $(1 - \pi)$  having mean zero and variance  $\sigma_j^2$ , which has a scaled inverse chi-square distribution; that is,  $\sigma_j^2 \sim \chi^{-2}(\nu_\alpha, S_\alpha^2)$ . Then, with  $\pi$  identical to 0, the marginal distribution of locus effects is a multivariate  $t$  distribution, and

BayesC is equal to RR-BLUP. In BayesC $\pi$ , the parameter  $\pi$  is an unknown quantity with a uniform prior, and  $\delta_j$  is used to indicate whether  $\alpha_j$  has a normal distribution ( $\delta_j=1$ ) or is zero ( $\delta_j=0$ ). Then, the parameter  $\pi$  can be sampled using a beta distribution with shape  $a = k - m + 1$  and shape  $b = m + 1$  with  $m = \sum_j \delta_j$ .

#### 4.1.5. LASSO

Tibshirani [20] developed a regression method, the least absolute shrinkage and selection operator (LASSO), that is a constrained version of ordinary least squares. It is somewhat indifferent to closely correlated markers and tends to pick one and ignore the others. The LASSO solution is the set of  $\alpha$  that satisfies

$$\min_{\alpha} \left\{ \sum (y_i - Z_i' \alpha)^2 + \lambda \sum_j |\alpha_j| \right\} \quad (2)$$

where  $\lambda \geq 0$ . Park and Casella [21] presented a Bayesian version of the LASSO (Bayesian LASSO) and suggested a Gibbs sampler for its implementation, with  $\alpha_j$  following a double-exponential prior:

$$p(\alpha_j | \lambda) = \frac{\lambda}{2} \exp(-\lambda |\alpha_j|) \quad (3)$$

However, Gibbs sampling often limits the speed of the regression. Using cyclical coordinate descent, Friedman et al. [22] developed fast algorithms for solving LASSO and the mixture of LASSO and ridge regression, named the elastic-net model.

#### 4.1.6. GBLUP

GBLUP is an efficient method that uses genome-wide markers to predict genetic and phenotypic values of candidates [7,23]. It assigns common variance to all loci and treats them as equally important in essence [24]. It uses a genomic relationship matrix (GRM) instead of the conventional pedigree-derived numerator relationship matrix  $A$ . The estimation of GEBV based on the GRM is more accurate [24–28]. VanRaden [24] suggested that the matrix  $G$  be established in the following form:

$$G = \frac{ZZ'}{2 \sum p_i (1 - p_i)} \quad (4)$$

where subtraction of  $P$  from  $M$  gives  $Z$ , and  $M$  is a genotypic matrix for markers coded as  $(-1, 0, 1)$ . The frequency of the second allele at locus  $i$  is  $p_i$ , and column  $i$  of  $P$  can be calculated as  $2(p_i - 0.5)$ . Then the breeding value  $a$  (equal to  $Z\alpha$ ) has a multivariate normal distribution with mean zero and covariance matrix  $G\sigma_a^2$ , where  $\sigma_a^2$  is the additive genetic variance. A normal distribution with mean zero and covariance matrix  $R\sigma_e^2$  is used for the residuals.

Yang et al. [27] used the following weighting scheme to calculate the matrix  $G$ :

$$G_{jk} = \frac{1}{N} \sum_i G_{ijk} = \begin{cases} \frac{1}{N} \sum_i \frac{(Z_{ij} - 2p_i)(Z_{ik} - 2p_i)}{2 \sum p_i (1 - p_i)}, & j \neq k \\ 1 + \frac{1}{N} \sum_i \frac{Z_{ij}^2 - (1 + 2p_i)Z_{ij} + 2p_i^2}{2 \sum p_i (1 - p_i)}, & j = k \end{cases} \quad (5)$$

where  $G_{ijk}$  is the relationship between the  $j$ th and  $k$ th individuals at locus  $i$ ;  $p_i$  is the allele frequency at SNP  $i$ , and  $Z_{ij}$  is an indicator variable that takes values of 0, 1, or 2 if the genotype of the  $j$ th individual at SNP  $i$  is  $bb$ ,  $Bb$  or  $BB$  (alleles are arbitrarily called  $b$  or  $B$ ), respectively.

Goddard et al. [28] used the relationship matrix  $A$  based on pedigree to calculate the matrix  $G$ :

$$G = A + b(G_m - A) \quad (6)$$

where  $b = \sigma_a^2 / \sigma_g^2$ ;  $\sigma_a^2$  is the total additive genetic variance, and  $\sigma_g^2$  is the additive genetic variance of each marker.  $G_m$  is calculated according to the algorithm proposed by VanRaden [24], and the encoding of  $Z$  is  $(0, 1, 2)$ .

Single step is another method for genetic evaluation which utilizes phenotypic, pedigree and genomic information [29]. This method constructs a relationship matrix that combines  $A$  and  $G$ , resulting in the following modified  $H$  [30]:

$$H = \begin{bmatrix} A_{11} + A_{12}A_{22}^{-1}(G - A_{22})A_{22}^{-1}A_{21} & A_{12}A_{22}^{-1}G \\ GA_{22}^{-1}A_{21} & G \end{bmatrix} \quad (7)$$

where the subscripts 1 and 2 of  $A$  denote ungenotyped and genotyped individuals, respectively.

#### 4.1.7. Reproducing kernel Hilbert space

In addition to the parameter methods introduced previously, there are other non-parametric and semi-parametric methods for GS, such as reproducing kernel Hilbert space (RKHS), which uses the Gauss kernel function to fit the following model [17]:

$$y = Xb + K_h \alpha + \epsilon \quad (8)$$

where  $\alpha$  has a multivariate normal distribution with mean zero and covariance matrix  $K_h \sigma_\alpha^2$ ;  $\epsilon \sim N(0, I_n \sigma^2)$ ;  $K_h$  is a kernel function that represents the correlation between individuals and is defined as.

$$K_h(x_i, x_j) = \exp(-hd_{ij}) \quad (9)$$

where  $d_{ij}$  is the squared Euclidean distance between individuals  $i$  and  $j$  calculated based on their genotypes, and the smoothing parameter  $h$  is defined as  $h = 2/d^*$  and  $d^*$  is the mean of  $d_{ij}$ . The model can be solved using a Gibbs sampler in a Bayesian framework, or using a mixed linear model.

#### 4.1.8. Support vector machine

Support vector machine (SVM) is a typical non-parametric method. It is a supervised learning method that can be used for classification and regression analysis. SVM follows the principle of structural risk minimization and takes into account the fitting and complexity of training samples. In recent years, it has been applied to genomic prediction [31]. When SVM is used for prediction analysis, a large data set with high dimension will lead to great computational complexity. The use of a kernel function greatly simplifies the inner product, so as to solve the curse of dimensionality. Thus, the selection of a kernel function is the key factor in SVM, and should reflect the distribution characteristics of the training samples. The commonly used kernel functions are linear kernel, Gauss RBF kernel, and polynomial kernel. The Gauss RBF kernel is widely adaptable and can be applied to any distribution of training samples with an appropriate



width parameter. Although it sometimes leads to overfitting problems, it is the most widely used kernel function. Some other machine learning methods, such as random forest (RF) and radial basis function neural network (RBFNN) can also be used in GS [32,33]. However, these methods are not widely used because they are complex and their predictive ability is similar to those of the common methods mentioned above.

#### 4.1.9. Characteristics of general methods

The assumption of GBLUP and RR-BLUP are that the effects of all loci have a common variance, making them more suitable for traits influenced by a large number of minor genes. Although Habier et al. [34] pointed out that GBLUP and RR-BLUP are fundamentally equivalent, GBLUP has two important characteristics. First, the dimensions of the model are reduced. The main computational difficulty is to evaluate the inverse of the  $n \times n$  matrix  $\mathbf{G} + \mathbf{R}(\frac{\sigma_g^2}{\sigma_a^2})$ . Second, using GBLUP, the GEBV can be estimated directly without thousands of iterations, which greatly increases efficiency.

However, in many cases, most markers in the genome have small or no effects, and a few markers have large effects, a situation that is not consistent with the hypothesis of GBLUP and RR-BLUP. Most Bayesian methods fit the actual situation, allowing different markers to have different effects and variances. The main difference between various Bayesian methods is that they have different prior distributions and then produce different degrees of shrinkage. When some loci have large effects on the trait, BayesB is the best choice because of its strongest shrinkage degree with a large value of  $\pi$ . In BayesC $\pi$ , the parameter  $\pi$  can be calculated from the posterior distribution based on experimental data and thus the shrinkage degree is determined, so it is more feasible than BayesB for real data analysis. The shrinkage degree of BayesA is weaker than BayesB and BayesC $\pi$ , and is thus more suitable for traits controlled by a moderate number of genes. Overall, Bayesian methods have their special priors that enable them to capture large-effect QTL and yield better predictive abilities. Studies have shown that the accuracy of Bayesian methods is slightly higher than that of GBLUP, especially for the scenario of distant relationships between training and testing sets [35,36].

In wheat, BayesB and the other selective shrinkage models, including BayesC $\pi$ , BayesA and Bayesian LASSO, were sensitive to the number of QTL, and the stronger the shrinkage degree, the more sensitive the models. Thus, accuracy decreases as the number of QTL increases [23]. Whereas, accuracy of GBLUP and RR-BLUP often stays nearly constant regardless of the number of QTL. They are more suitable for traits controlled by a large number of minor genes because of the strong robustness. BayesA was found to be widely adaptable because of its moderate shrinkage degree [23]. However, the estimation of Bayesian methods is often time-consuming, restricting their application. A better alternative is the fast algorithm for LASSO [22], because it strikes a balance between selective shrinkage and computational efficiency. Although previous studies [37] in Nellore cattle showed that Bayesian regression models were more accurate than GBLUP, Xu et al. [38] pointed out that GBLUP was superior to Bayesian methods in the prediction of six yield-related traits for maize, and prediction of grain yield in wheat also

achieved similar results [23]. A reasonable explanation is that most quantitative traits in crops are complex and accuracy of various methods is associated with specific traits of a particular group.

In GS, nonparametric methods such as machine learning have also been successfully applied. A comparison of various GS methods in a mouse population showed good performance of RKHS for prediction [17]. In wheat, studies of days to heading and grain yield demonstrated a consistent superiority of RKHS and RBFNN over the Bayesian LASSO, Bayesian ridge regression, BayesA, and BayesB models [39]. In chickpea, RR-BLUP, BayesC $\pi$ , BayesB, Bayesian LASSO and RF resulted in similar accuracy for traits of interest [40]. Although the most suitable method is usually case-dependent, Gonzalez-Recio et al. [41] suggested that SVM and RF are good at solving classification problems, whereas RKHS is suitable for regression problems.

#### 4.2. Prediction methods for non-additive genetic effects

Depending on the genetic effects to be estimated, there are two different strategies in GS. One focuses on estimating breeding value (BV), the additive effects that can be directly transmitted from parents to offspring. Non-additive effects such as dominance and epistasis are associated with specific genotypes and cannot be inherited consistently. When variance components are estimated, non-additive effects are often incorporated into the random environmental effects and considered to be noise. The other strategy focuses on both additive and non-additive effects and can help us explore heterosis. Many studies have shown that incorporating non-additive effects could be very beneficial in some populations [42]. If non-additive effects are significant, their neglect will lead to bias in genetic estimation [43].

Because non-additive genetic effects may make an important contribution to total genetic variation of complex traits, it is necessary to consider them in GS. Dominance estimation is essential for vegetatively propagated species [44] and crossed populations, where including both additive effects and dominance might contribute to prediction accuracy and serve as a guide for choosing mating pairs [45]. Balestre et al. [46] detected high dominance effects in predicting grain yield of maize single crosses using the additive and dominance model based on BLUP. Additionally, the interactions between genes are well documented, and the importance of epistasis is still an area of active research. In near-isogenic rice lines, epistasis was detected among three major flowering-time genes [47]. Mao et al. [48] reported that digenic epistasis accounted for important genetic bases of transgressive segregation for thousand-grain weight in a RIL rice population. Dudley and Johnson [49] found that epistatic effects were important in determination of oil, protein, and starch contents of corn. They concluded that epistasis was an important contributor to the long-term response to selection of these quantitative traits. In rapeseed, Wuerschum et al. [50] observed a high contribution of epistasis to the genotypic variance of complex quantitative traits. Hu et al. [51] predicted the genomic value of somatic embryo number for soybean RIL lines, and cross-validation showed that prediction ability was greatly increased when epistatic effects were included in the

model. Prediction of genetic values for 280 wheat accessions also showed that prediction accuracy could be substantially improved by the inclusion of two-locus epistatic effects [42]. In summary, the potential advantage of incorporating non-additive effects should receive attention in GS research for crop breeding.

#### 4.3. Prediction methods on multiple traits

Most previous empirical studies with GS focused on within-environment predictions based on single-environment (SE) models. These models typically target a single phenotypic trait, disregarding the reality that multiple traits or phenotypes in multiple environments are most likely associated. Simultaneously modeling multiple traits or multiple environments can make use of genetically and environmentally correlated information and improve the accuracy of prediction [52–54]. Taking the multi-trait (MT) scenario as an example (for  $s$  traits), the classical multivariate (MV) model can be described in reduced matrix form as.

$$y = Xb + Z_a a + \varepsilon \quad (10)$$

where  $y = [y_1^T, y_2^T, \dots, y_s^T]^T$ ;  $b = [b_1^T, b_2^T, \dots, b_s^T]^T$ ;  $a = [a_1^T, a_2^T, \dots, a_s^T]^T$  and  $\varepsilon = [\varepsilon_1^T, \varepsilon_2^T, \dots, \varepsilon_s^T]^T$ . The non-genetic effects  $b$  are treated as fixed effects. The additive effects  $a$  and residuals  $\varepsilon$  are treated as random effects following multivariate normal distributions; specifically:  $a \sim N(0, G_{a0} \otimes G_{a_s})$  and  $\varepsilon \sim N(0, R_{\varepsilon} \otimes I_m \sigma_{\varepsilon_s}^2)$ , where  $G$  is the genomic relationship matrix,  $\otimes$  denotes the Kronecker product of matrices [55],  $m$  is the number of phenotypic observations, and  $I_m$  is an  $m \times m$  identity matrix.  $X$  and  $Z_a$  are incidence matrices for the fixed effects and random additive effects, respectively. The additive genetic covariance matrix  $G_{a0}$  and  $R_{\varepsilon}$  can be written as:

$$G_{a0} = \begin{bmatrix} \sigma_{a_1}^2 & \rho_{a12}\sigma_{a_1}\sigma_{a_2} & \dots & \rho_{a1s}\sigma_{a_1}\sigma_{a_s} \\ \rho_{a21}\sigma_{a_2}\sigma_{a_1} & \sigma_{a_2}^2 & \dots & \rho_{a2s}\sigma_{a_2}\sigma_{a_s} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{as1}\sigma_{a_s}\sigma_{a_1} & \rho_{as2}\sigma_{a_s}\sigma_{a_2} & \dots & \sigma_{a_s}^2 \end{bmatrix} \quad (11)$$

$$R_{\varepsilon} = \begin{bmatrix} \sigma_{\varepsilon_1}^2 & \rho_{12}\sigma_{\varepsilon_1}\sigma_{\varepsilon_2} & \dots & \rho_{1s}\sigma_{\varepsilon_1}\sigma_{\varepsilon_s} \\ \rho_{21}\sigma_{\varepsilon_2}\sigma_{\varepsilon_1} & \sigma_{\varepsilon_2}^2 & \dots & \rho_{2s}\sigma_{\varepsilon_2}\sigma_{\varepsilon_s} \\ \vdots & \vdots & \ddots & \vdots \\ \rho_{s1}\sigma_{\varepsilon_s}\sigma_{\varepsilon_1} & \rho_{s2}\sigma_{\varepsilon_s}\sigma_{\varepsilon_2} & \dots & \sigma_{\varepsilon_s}^2 \end{bmatrix} \quad (12)$$

where  $\sigma_{a_i}^2$  and  $\sigma_{\varepsilon_i}^2$  are the additive and residual variance of lines with trait  $i$  ( $i = 1, 2, \dots, s$ ), respectively.  $\rho_{aij}$  and  $\rho_{ij}$  are the additive and residual correlations of lines between traits  $i$  and  $j$ , respectively.

Multiple traits are most likely associated owing to pleiotropy and shared biological basis [56]. Actual crop breeding often targets multiple correlated traits, and joint analysis taking into consideration the correlation may result in more accurate prediction. GS has been developed from single-trait (ST) to MT models. Simultaneously modeling multiple quantitative traits results in better predictive power than individually targeting traits [57,58]. Using single-step methods, Tsuruta et al. [59] found that the MT model is more accurate than the ST model. Previous studies of GS for crops [52] showed that MT models had higher accuracy than ST models when phenotypes are not available for all individuals and traits. With simulated data, Guo et al. [60] found that the MT

model performed better than the ST model for traits with low heritability for a limited number of records. For rice hybrids, when phenotypes for a focal trait were unknown and that for the other traits were observed (unbalanced scenarios), a study [15] showed that the average predictive ability of the MT model with two traits was 6.4% higher than that of a univariate (UV) model, and the gain of the MT model with eight traits over the UV model was 26.7%. When predicting the individuals without any phenotypes (balanced scenarios), the benefit of prediction is far lower than that in unbalanced scenarios [54,61], which restricts the wide application of MT prediction. However, unbalanced MT prediction can bring new strategies to crop breeding. Most economic traits of crops are quantitative traits and controlled by multiple genes. They are readily affected by environment and their phenotypes are difficult to measure accurately. Some traits are expensive to identify, such as the crop yield and trace element levels in grain, and others are difficult to identify early in development, such as cold resistance in the seedling stage of rice, wheat and corn. MT prediction using auxiliary traits that are easier or cheaper to record represents an effective strategy that improves accuracy and saves phenotyping costs.

In crop breeding, multi-environment (ME) trials play a fundamental role in assessing the performance of genotypes across different environmental conditions, and genotype-by-environment ( $G \times E$ ) interaction has always been a concern. ME models allow sharing of information across environments, and previous studies [55,62,63] have illustrated that the use of genetic and residual covariance across correlated environments might improve cross-environment accuracy in ME models, especially for varieties evaluated in some environments but not in others. In wheat and maize, a recent study [64] confirmed that ME models were predictively uniformly superior to SE models.

SI is also a common method for breeding selection of multiple traits. It is generally more efficient than selection based on independent culling levels or tandem selection, which can improve several quantitative traits simultaneously. Based on classical SI, Kempthorne et al. [65] proposed a restricted SI method that is designed for a situation when breeders want to improve some traits but simultaneously hold other traits at their average levels. Recently, Cerón-Rojas et al. [66] proposed a more efficient predetermined proportional gains eigen-SI method. SI methods have been widely used in crops such as alfalfa, oat, maize, and soybean [67–70]. The rapid development of GS has brought new perspectives for SI. Dekkers [71] incorporated marker information in standard software for SI predictions of response and rates of inbreeding. Cerón-Rojas et al. [72] developed the theory of a genomic selection index and applied it to simulated and real data sets of maize. Schulthess et al. [54] applied GS to predict different SIs for grain yield and protein content improvement in rye. Lyra et al. [73] evaluated predictions of four SIs that could contribute to the selection of tropical maize hybrids under contrasting nitrogen conditions, and the results showed that the method was effective.

## 5. Factors affecting the accuracy of genomic selection

The accuracy of GS is affected by many factors, such as sample size, genetic relationship, marker density, heritability,

and LD between markers and QTL. Usually, a large training set can improve the accuracy of prediction, because the larger the training set, the more samples with phenotypic and genotypic information can be used, and thus the accuracy of estimation of genetic effects can be improved accordingly [5], especially for traits with low heritability [74].

The level of relatedness between genotypes in the training set and the testing set also has a strong impact on prediction accuracy [26,75]. In an interconnected biparental maize population, Riedelsheimer et al. [4] investigated the influence of the composition of the training set on genomic prediction accuracy. Prediction accuracy declined by 42% if full-sib were replaced by half-sib DH lines, but statistically significantly better results could be achieved if half-sib DH lines were available from both parents instead of only one parent of the validation population.

GS uses markers covering the whole genome so that all genetic variance can be explained by the markers. It is assumed that there are always some markers in LD with any QTL. Sufficiently high marker density guarantees near-perfect LD between at least one marker and each QTL, leading to higher prediction accuracy. Although in theory, the more markers there are, the better the prediction, accuracy is difficult to improve meaningfully when the density reaches a certain degree. Su et al. [76] investigated genomic prediction using medium-density and high-density marker panels. When marker density increased from 54 K to 777 K, prediction accuracy increased only 0.5–1.0%. For grain quality traits in biparental wheat populations, studies showed that accuracy for GS reached a plateau at low marker densities (128–256) [77]. A similar tendency is found in other studies [15,27] when the marker density reaches a certain level, it is no longer beneficial to genomic prediction. In short, with an increase in marker density, the accuracy of GS is not proportional to the cost.

When the density is fixed, the length of haplotypes may also affect accuracy. In two studies [78,79], haplotypes made up of 10 markers yielded the highest estimation accuracy for breeding values. The degree of LD between markers and QTL may also affect accuracy. If GS is applied for many generations, the effect of the markers does not change but the proportion of the genetic variance explained by them declines, and accuracy may decrease rapidly in subsequent generations after estimation owing to decay of genetic relationships [34]. Meuwissen et al. [18] demonstrated that the accuracy of GS in the first two generations decreased rapidly, while accuracy in the succeeding generations decreased slowly. With the passage of generations, accuracy for traits with high heritability decreases more slowly [79]. To ensure high accuracy in GS, it is best to re-estimate the marker or haplotype effects every 2–3 generations. However, Zhong et al. [80] pointed out that predictions relied primarily on marker information to model genetic relationships between the training and testing sets, rather than on markers capturing QTL effects by association, suggesting that relationship information was more valuable than LD information.

The accuracy of GS also depends on the genetic architecture of traits, such as heritability and distribution of causal genes. Heritability is positively related to prediction accuracy. Using the same GS method, prediction accuracy of a high-heritability trait (such as thousand-grain weight) is often higher than that of a low-heritability trait (such as grain yield).

$G \times E$  interaction is another important factor. For example, accuracies for the same trait in various environments differ greatly because of interaction. With a winter wheat data set consisting of 2437 genotypes tested for grain yield in 44 environments, Heslot et al. [63] randomly split it into two sets of 22 environments, finding that the accuracy in predicting genotype performance in unobserved environments increased by 11.1% on average when weather data were available.

## 6. Perspective

Using substantially different means, various GS methods deal with the issues of increasing dimensionality and computational complexity, and thereby capture different aspects of the association between genotype and phenotype. For this reason, the performances of different methods depend on the genetic architecture underlying the specific trait. For crop improvement, most quantitative traits are influenced by polygenes. Among the general GS methods, GBLUP is recommended for breeding practice because of its robustness, and computational efficiency. Some machine learning algorithms, such as RKHS, have also been successful in GS, and likely represent avenues for future research.

A synthesis of various models may be superior to any single model. Jannink et al. [81] reported that a simple mean across all methods unexpectedly performed best. The classical elastic-net model is just a mixture of two different methods [22]. Using genetic architecture information contained in the dataset is another effective strategy. Zhang et al. [7] proposed a method, termed BLUP|GA, that can make full use of a trait-specific covariance matrix, and showed that BLUP|GA outperformed GBLUP, BayesA, and BayesB in most cases.

Although the interactions between genes are well documented, and the role of epistasis in expression of complex traits has been widely recognized, analysis of genome-wide loci and their interaction involves many variables, far more than the number of observed samples. Internationally, some statistical methods such as LASSO and GBLUP have been applied to these big-data analyses [13,22]. Nonetheless, a large number of markers with epistasis pose both statistical and computational challenges to prediction. For this reason, most studies that attempt to identify the genetic basis of complex traits ignore epistatic effects. However, in breeding practice, the simplification of a model may mask some genetic variance, resulting in biased and unreliable analysis in GS. How to accurately define the interactions that occur among many genes and to efficiently estimate the effects remain subjects of GS research.

Joint analysis of multiple traits can improve the accuracy of prediction with highly correlated traits, particularly for some low-heritability traits [52,82]. A previous study has shown that MT models with more auxiliary traits might greatly improve predictive ability [15]. Moreover, ME models can make use of environmental covariates and consider  $G \times E$  as a lack of genetic correlation between environments, so as to obtain higher accuracy than the ST prediction. With the rapid development of phenotypic omics and environmental omics, breeders can acquire more phenotypic and environmental information at lower cost, providing great opportunity for the research of MT and ME prediction.



Not only genomic, but also metabolomic and transcriptomic data have received attention for phenotypic prediction. Riedelsheimer et al. [4] exploited genomic and metabolic information to predict complex traits in maize hybrid test crosses. Xu et al. [38] compared the predictive abilities of three different types of omic data, including genomic, transcriptomic, and metabolomic data, for 339 maize inbred lines. In rice, studies [83] showed that the predictive power of biomass heterosis for hybrids was significantly improved with parental metabolic measurements. In short, as the results of multiple genes and their products, the metabolome and transcriptome have proven to be useful in prediction. Integrating multiple omic data is expected to be an important method for the study of GS in crop breeding.

## Acknowledgments

This work was supported by grants from the National High Technology Research and Development Program of China (2014AA10A601-5), the National Key Research and Development Program of China (2016YFD0100303), the National Natural Science Foundation of China (91535103), the Natural Science Foundations of Jiangsu Province (BK20150010), the Natural Science Foundation of the Jiangsu Higher Education Institutions (14KJA210005), the Open Research Fund of State Key Laboratory of Hybrid Rice (Wuhan University) (KF201701), the Science and Technology Innovation Fund Project in Yangzhou University (2016CXJ021), the Priority Academic Program Development of Jiangsu Higher Education Institutions and the Innovative Research Team of Universities in Jiangsu Province.

## REFERENCES

- [1] L.N. Hazel, J.L. Lush, The efficiency of three methods of selection, *J. Hered.* 33 (1942) 393–399.
- [2] C.R. Henderson, Best linear unbiased estimation and prediction under a selection model, *Biometrics* 31 (1975) 423–447.
- [3] R. Lande, R. Thompson, Efficiency of marker-assisted selection in the improvement of quantitative traits, *Genetics* 124 (1990) 743–756.
- [4] C. Riedelsheimer, A. Czedik-Eysenberg, C. Grieder, J. Lisec, F. Technow, R. Sulpice, T. Altmann, M. Stitt, L. Willmitzer, A.E. Melchinger, Genomic and metabolic prediction of complex heterotic traits in hybrid maize, *Nat. Genet.* 44 (2012) 217–220.
- [5] M. Goddard, Genomic selection: prediction of accuracy and maximisation of long term response, *Genetica* 136 (2009) 245–257.
- [6] L. Chen, C. Li, M. Sargolzaei, F. Schenkel, Impact of genotype imputation on the performance of GBLUP and Bayesian methods for genomic prediction, *PLoS One* 9 (2014), e101544.
- [7] Z. Zhang, M. Erbe, J. He, U. Ober, N. Gao, H. Zhang, H. Simianer, J. Li, Accuracy of whole-genome prediction using a genetic architecture-enhanced variance-covariance matrix, *G3-Genes Genomes Genet.* 5 (2015) 615–627.
- [8] E.L. Heffner, A.J. Lorenz, J.L. Jannink, M.E. Sorrells, Plant breeding with genomic selection: gain per unit time and cost, *Crop Sci.* 50 (2010) 1681–1690.
- [9] R. Bernardo, J. Yu, Prospects for genomewide selection for quantitative traits in maize, *Crop Sci.* 47 (2007) 1082–1090.
- [10] T. Albrecht, V. Wimmer, H. Auinger, M. Erbe, C. Knaak, M. Ouzunova, H. Simianer, C.C. Schön, Genome-based prediction of testcross values in maize, *Theor. Appl. Genet.* 123 (2011) 339–350.
- [11] D.C. Kadam, S.M. Potts, M.O. Bohn, A.E. Lipka, A.J. Lorenz, Genomic prediction of single crosses in the early stages of a maize hybrid breeding pipeline, *G3-Genes Genomes Genet.* 6 (2016) 3443–3453.
- [12] U. Beukert, Z. Li, G. Liu, Y. Zhao, N. Ramachandra, V. Mirdita, F. Pita, K. Pillen, J.C. Reif, Genome-based identification of heterotic patterns in rice, *Rice* 10 (2017) 22.
- [13] S.Z. Xu, D. Zhu, Q.F. Zhang, Predicting hybrid performance in rice using genomic best linear unbiased prediction, *Proc. Natl. Acad. Sci. U. S. A.* 111 (2014) 12456–12461.
- [14] H.D. Daetwyler, U.K. Bansal, H.S. Bariana, M.J. Hayden, B.J. Hayes, Genomic prediction for rust resistance in diverse wheat landraces, *Theor. Appl. Genet.* 127 (2014) 1795–1803.
- [15] X. Wang, L.Z. Li, Z.F. Yang, X.F. Zheng, B. Yu, C.W. Xu, Z.L. Hu, Predicting rice hybrid performance using univariate and multivariate GBLUP models based on North Carolina mating design II, *Heredity* 118 (2017) 302–310.
- [16] J.Y. Li, J. Wang, R.S. Zeigler, The 3,000 rice genomes project: new opportunities and challenges for future rice research, *Gigascience* 3 (2014) 8.
- [17] H.H. Neves, R. Carvalheiro, S.A. Queiroz, A comparison of statistical methods for genomic selection in a mice population, *BMC Genet.* 13 (2012) 100.
- [18] T.H. Meuwissen, B.J. Hayes, M.E. Goddard, Prediction of total genetic value using genome-wide dense marker maps, *Genetics* 157 (2001) 1819–1829.
- [19] S.Z. Xu, Estimating polygenic effects using markers of the entire genome, *Genetics* 163 (2003) 789–801.
- [20] R. Tibshirani, Regression shrinkage and selection via the lasso, *J. R. Stat. Soc. B* (1996) 267–288.
- [21] T. Park, G. Casella, The bayesian lasso, *J. Am. Stat. Assoc.* 103 (2008) 681–686.
- [22] J. Friedman, T. Hastie, R. Tibshirani, Regularization paths for generalized linear models via coordinate descent, *J. Stat. Softw.* 33 (2010) 1–22.
- [23] X. Wang, Z.F. Yang, C.W. Xu, A comparison of genomic selection methods for breeding value prediction, *Sci. Bull.* 60 (2015) 925–935.
- [24] P. VanRaden, Efficient methods to compute genomic predictions, *J. Dairy Sci.* 91 (2008) 4414–4423.
- [25] B.J. Hayes, P.M. Visscher, M.E. Goddard, Increased accuracy of artificial selection by using the realized relationship matrix, *Genet. Res.* 91 (2009) 47–60.
- [26] D. Habier, J. Tetens, F.-R. Seefried, P. Lichtner, G. Thaller, The impact of genetic relationship information on genomic breeding values in German Holstein cattle, *Genet. Sel. Evol.* 42 (2010) 5.
- [27] J. Yang, B. Benyamin, B.P. McEvoy, S. Gordon, A.K. Henders, D. R. Nyholt, P.A. Madden, A.C. Heath, N.G. Martin, G.W. Montgomery, Common SNPs explain a large proportion of the heritability for human height, *Nat. Genet.* 42 (2010) 565–569.
- [28] M.E. Goddard, B.J. Hayes, T.H.E. Meuwissen, Using the genomic relationship matrix to predict the accuracy of genomic selection, *J. Anim. Breed. Genet.* 128 (2011) 409–421.
- [29] O.F. Christensen, P. Madsen, B. Nielsen, T. Ostensen, G. Su, Single-step methods for genomic evaluation in pigs, *Animal* 6 (2012) 1565–1571.
- [30] A. Legarra, I. Aguilar, I. Misztal, A relationship matrix including full pedigree and genomic information, *J. Dairy Sci.* 92 (2009) 4656–4663.
- [31] S. Maenhout, B. De Baets, G. Haesaert, E. van Bockstaele, Support vector machine regression for the prediction of maize hybrid performance, *Theor. Appl. Genet.* 115 (2007) 1003–1013.
- [32] N. Heslot, H.P. Yang, M.E. Sorrells, J.L. Jannink, Genomic selection in plant breeding: a comparison of models, *Crop Sci.* 52 (2012) 146–160.

- [33] J. Crossa, P. Perez-Rodriguez, J. Cuevas, O. Montesinos-Lopez, D. Jarquin, G. de los Campos, J. Burgueno, J.M. Gonzalez-Camacho, S. Perez-Elizalde, Y. Beyene, S. Dreisigacker, R. Singh, X.C. Zhang, M. Gowda, M. Roorkiwal, J. Rutkoski, R.K. Varshney, Genomic selection in plant breeding: methods, models, and perspectives, *Trends Plant Sci.* 22 (2017) 961–975.
- [34] D. Habier, R. Fernando, J. Dekkers, The impact of genetic relationship information on genome-assisted breeding values, *Genetics* 177 (2007) 2389–2397.
- [35] H. Gao, G. Su, L. Janss, Y. Zhang, M.S. Lund, Model comparison on genomic predictions using high-density markers for different groups of bulls in the Nordic Holstein population, *J. Dairy Sci.* 96 (2013) 4678–4687.
- [36] X. Wu, M.S. Lund, D. Sun, Q. Zhang, G. Su, Impact of relationships between test and training animals and among training animals on reliability of genomic prediction, *J. Anim. Breed. Genet.* 132 (2015) 366–375.
- [37] H.H.R. Neves, R. Carvalho, A.M. Perez, Y.T. O'Brien, A.S. Utsunomiya, F.S. do Carmo, J. Schenkel, J.C. Soelkner, C.P. McEwan, J.B. Van Tassell, M.V.G.B. da Cole, S.A. Silva, T.S. Queiroz, J.F. Garcia Sonstegard, Accuracy of genomic predictions in *Bos indicus* (Nelore) cattle, *Genet. Sel. Evol.* 46 (2014) 17.
- [38] Y. Xu, C.W. Xu, S.Z. Xu, Prediction and association mapping of agronomic traits in maize using multiple omic data, *Heredity* 119 (2017) 174–184.
- [39] P. Perez-Rodriguez, D. Gianola, J.M. Gonzalez-Camacho, J. Crossa, Y. Manes, S. Dreisigacker, Comparison between linear and non-parametric regression models for genome-enabled prediction in wheat, *G3-Genes Genomes Genet.* 2 (2012) 1595–1605.
- [40] M. Roorkiwal, A. Rathore, R.R. Das, M.K. Singh, A. Jain, S. Srinivasan, P.M. Gaur, B. Chellapilla, S. Tripathi, Y. Li, Genome-enabled prediction models for yield related traits in chickpea, *Front. Plant Sci.* 7 (2016) 1666.
- [41] O. Gonzalez-Recio, G.J.M. Rosa, D. Gianola, Machine learning methods and predictive ability metrics for genome-wide prediction of complex traits, *Livest. Sci.* 166 (2014) 217–231.
- [42] D. Wang, I.S. El-Basyoni, P.S. Baenziger, J. Crossa, K.M. Eskridge, I. Dweikat, Prediction of genetic values of quantitative traits with epistatic effects in plant breeding populations, *Heredity* 109 (2012) 313–319.
- [43] I. Misztal, T.J. Lawlor, R.L. Fernando, Dominance models with method R for stature of Holsteins, *J. Dairy Sci.* 80 (1997) 975–978.
- [44] M. Denis, J.M. Bouvet, Efficiency of genomic selection with models including dominance effect in the context of eucalyptus breeding, *Tree Genet. Genomes* 9 (2013) 37–51.
- [45] R. Wellmann, J. Bennewitz, Bayesian models with dominance effects for genomic evaluation of quantitative traits, *Genet. Res.* 94 (2012) 21–37.
- [46] M. Balestre, R.G. Von Pinho, J.C. Souza, Prediction of maize single-cross performance by mixed linear models with microsatellite marker information, *Genet. Mol. Res.* 9 (2010) 1054–1068.
- [47] N. Uwatoko, A. Onishi, Y. Ikeda, M. Kontani, A. Sasaki, K. Matsubara, Y. Itoh, Y. Sano, Epistasis among the three major flowering time genes in rice: coordinate changes of photoperiod sensitivity, basic vegetative growth and optimum photoperiod, *Euphytica* 163 (2008) 167–175.
- [48] D. Mao, T. Liu, C. Xu, X. Li, Y. Xing, Epistasis and complementary gene action adequately account for the genetic bases of transgressive segregation of kilo-grain weight in rice, *Euphytica* 180 (2011) 261–271.
- [49] J.W. Dudley, G.R. Johnson, Epistatic models improve prediction of performance in corn, *Crop Sci.* 49 (2009) 1533.
- [50] T. Würschum, H.P. Maurer, F. Dreyer, J.C. Reif, Effect of inter- and intragenic epistasis on the heritability of oil content in rapeseed (*Brassica napus* L.), *Theor. Appl. Genet.* 126 (2013) 435–441.
- [51] Z. Hu, Y. Li, X. Song, Y. Han, X. Cai, S. Xu, W. Li, Genomic value prediction for quantitative traits under the epistatic model, *BMC Genet.* 12 (2011) 15.
- [52] Y. Jia, J.L. Jannink, Multiple-trait genomic selection methods increase genetic value prediction accuracy, *Genetics* 192 (2012) 1513–1522.
- [53] M. Lopez-Cruz, J. Crossa, D. Bonnett, S. Dreisigacker, J. Poland, J.L. Jannink, R.P. Singh, E. Autrique, G. de los Campos, Increased prediction accuracy in wheat breeding trials using a marker  $\times$  environment interaction genomic selection model, *G3-Genes Genomes Genet.* 5 (2015) 569–582.
- [54] A.W. Schulthess, Y. Wang, T. Miedaner, P. Wilde, J.C. Reif, Y. Zhao, Multiple-trait and selection indices-genomic predictions for grain yield and protein content in rye for feeding purposes, *Theor. Appl. Genet.* 129 (2016) 273–287.
- [55] Z. Guo, D.M. Tucker, D. Wang, C.J. Basten, E. Ersoz, W.H. Briggs, J. Lu, M. Li, G. Gay, Accuracy of across-environment genome-wide prediction in maize nested association mapping populations, *G3-Genes Genomes Genet.* 3 (2013) 263–272.
- [56] M. Scutari, P. Howell, D.J. Balding, I. Mackay, Multiple quantitative trait analysis using Bayesian networks, *Genetics* 198 (2014) 129–137.
- [57] C. Henderson, R. Quaas, Multiple trait evaluation using relatives' records, *J. Anim. Sci.* 43 (1976) 1188–1197.
- [58] T. Hayashi, H. Iwata, A Bayesian method and its variational approximation for prediction of genomic breeding values in multiple traits, *BMC Bioinf.* 14 (2013) 34.
- [59] S. Tsuruta, I. Misztal, I. Aguilar, T.J. Lawlor, Multiple-trait genomic evaluation of linear type traits using genomic and phenotypic data in US Holsteins, *J. Dairy Sci.* 94 (2011) 4198–4204.
- [60] G. Guo, F. Zhao, Y. Wang, Y. Zhang, L. Du, G. Su, Comparison of single-trait and multiple-trait genomic prediction models, *BMC Genet.* 15 (2014) 30.
- [61] Y. Bao, J.E. Kurlle, G. Anderson, N.D. Young, Association mapping and genomic prediction for resistance to sudden death syndrome in early maturing soybean germplasm, *Mol. Breed.* 35 (2015) 128.
- [62] J. Burgueno, G. de los Campos, K. Weigel, J. Crossa, Genomic prediction of breeding values when modeling genotype  $\times$  environment interaction using pedigree and dense molecular markers, *Crop Sci.* 52 (2012) 707–719.
- [63] N. Heslot, D. Akdemir, M.E. Sorrells, J.L. Jannink, Integrating environmental covariates and crop modeling into the genomic selection framework to predict genotype by environment interactions, *Theor. Appl. Genet.* 127 (2014) 463–480.
- [64] J. Cuevas, J. Crossa, O.A. Montesinos-Lopez, J. Burgueno, P. Perez-Rodriguez, G. de los Campos, Bayesian genomic prediction with genotype  $\times$  environment interaction kernel models, *G3-Genes Genomes Genet.* 7 (2017) 41–53.
- [65] O. Kempthorne, A.W. Nordskog, Restricted selection indices, *Biometrics* 15 (1959) 10–19.
- [66] J.J. Cerón-Rojas, J. Crossa, F.H. Toledo, J. Sahagún-Castellanos, A predetermined proportional gains eigen selection index method, *Crop Sci.* 56 (2016) 2436–2447.
- [67] J. Elgin, R. Hill, K. Zeiders, Comparison of four methods of multiple trait selection for five traits in alfalfa, *Crop Sci.* 10 (1970) 190–193.
- [68] D.J. Dolan, D.D. Stuthman, F.L. Kolb, A.D. Hewings, Multiple trait selection in a recurrent selection population in oat (*Avena sativa* L.), *Crop Sci.* 36 (1996) 1207–1211.
- [69] K. Suwantaradon, S. Eberhart, J. Mock, J. Owens, W. Guthrie, Index selection for several agronomic traits in the BSSS2 maize population, *Crop Sci.* 15 (1975) 827–833.
- [70] C.C. Holbrook, J.W. Burton, T.E. Carter, Evaluation of recurrent restricted index selection for increasing yield while holding seed protein constant in soybean, *Crop Sci.* 29 (1989) 324–329.
- [71] J.C.M. Dekkers, Prediction of response to marker-assisted and genomic selection using selection index theory, *J. Anim. Breed. Genet.* 124 (2007) 331–341.

- [72] J.J. Cerón-Rojas, J. Crossa, V.N. Arief, K. Basford, J. Rutkoski, D. Jarquin, G. Alvarado, Y. Beyene, K. Semagn, I. DeLacy, A genomic selection index applied to simulated and real data, *G3-Genes Genomes Genet.* 5 (2015) 2155–2164.
- [73] D.H. Lyra, L. de Freitas Mendonça, G. Galli, F.C. Alves, Í.S.C. Granato, R. Fritsche-Neto, Multi-trait genomic prediction for nitrogen response indices in tropical maize hybrids, *Mol. Breed.* 37 (2017) 80.
- [74] H.M. Nielsen, A.K. Sonesson, H. Yazdi, T.H.E. Meuwissen, Comparison of accuracy of genome-wide and BLUP breeding value estimates in sib based aquaculture breeding schemes, *Aquaculture* 289 (2009) 259–264.
- [75] S.A. Clark, J.M. Hickey, H.D. Daetwyler, J.H.J. van der Werf, The importance of information on relatives for the prediction of genomic breeding values and the implications for the makeup of reference data sets in livestock breeding schemes, *Genet. Sel. Evol.* 44 (2012) 4.
- [76] G. Su, R.F. Brondum, P. Ma, B. Guldbandsen, G.R. Aamand, M.S. Lund, Comparison of genomic predictions using medium-density (similar to 54,000) and high-density (similar to 777,000) single nucleotide polymorphism marker panels in Nordic Holstein and Red Dairy Cattle populations, *J. Dairy Sci.* 95 (2012) 4657–4665.
- [77] E.L. Heffner, J.L. Jannink, H. Iwata, E. Souza, M.E. Sorrells, Genomic selection accuracy for grain quality traits in biparental wheat populations, *Crop Sci.* 51 (2011) 2597.
- [78] M.P.L. Calus, T.H.E. Meuwissen, A.P.W. de Roos, R.F. Veerkamp, Accuracy of genomic selection using different methods to define haplotypes, *Genetics* 178 (2008) 553–561.
- [79] T.M. Villumsen, L. Janss, M.S. Lund, The importance of haplotype length and heritability using genomic selection in dairy cattle, *J. Anim. Breed. Genet.* 126 (2009) 3–13.
- [80] S.Q. Zhong, J.C.M. Dekkers, R.L. Fernando, J.L. Jannink, Factors affecting accuracy from genomic selection in populations derived from multiple inbred lines: a barley case study, *Genetics* 182 (2009) 355–364.
- [81] J.L. Jannink, A.J. Lorenz, H. Iwata, Genomic selection in plant breeding: from theory to practice, *Brief. Funct. Genomics* 9 (2010) 166–177.
- [82] N.A. Alimi, M.C.A.M. Bink, J.A. Dieleman, J.J. Magán, A.M. Wubs, A. Palloix, F.A.V. Eeuwijk, Multi-trait and multi-environment QTL analyses of yield and a set of physiological traits in pepper, *Theor. Appl. Genet.* 126 (2013) 2597–2625.
- [83] T. Gartner, M. Steinfath, S. Andorf, J. Lisec, R.C. Meyer, T. Altmann, L. Willmitzer, J. Selbig, Improved heterosis prediction by combining information on DNA- and metabolic markers, *PLoS One* 4 (2009), e5220.