

METHODOLOGY ARTICLE

Open Access



# Efficient strategies for leave-one-out cross validation for genomic best linear unbiased prediction

Hao Cheng<sup>1,2</sup>, Dorian J. Garrick<sup>1,3</sup> and Rohan L. Fernando<sup>1\*</sup>

## Abstract

**Background:** A random multiple-regression model that simultaneously fit all allele substitution effects for additive markers or haplotypes as uncorrelated random effects was proposed for Best Linear Unbiased Prediction, using whole-genome data. Leave-one-out cross validation can be used to quantify the predictive ability of a statistical model.

**Methods:** Naive application of Leave-one-out cross validation is computationally intensive because the training and validation analyses need to be repeated  $n$  times, once for each observation. Efficient Leave-one-out cross validation strategies are presented here, requiring little more effort than a single analysis.

**Results:** Efficient Leave-one-out cross validation strategies is 786 times faster than the naive application for a simulated dataset with 1,000 observations and 10,000 markers and 99 times faster with 1,000 observations and 100 markers. These efficiencies relative to the naive approach using the same model will increase with increases in the number of observations.

**Conclusions:** Efficient Leave-one-out cross validation strategies are presented here, requiring little more effort than a single analysis.

**Keywords:** Leave-one-out cross validation, GBLUP

## Background

A random multiple-regression model that simultaneously fit all allele substitution effects for additive markers or haplotypes as uncorrelated random effects was proposed for Best Linear Unbiased Prediction (BLUP) [1], using whole-genome data. Breeding values are defined as the sum of the effects of all the markers or haplotypes, and their estimates are widely used for prediction of the merit of selection candidates. Estimates of marker or haplotype effects are used to predict breeding values of individuals that were not present in a previous analysis commonly referred to as training. An alternative earlier published approach to use marker or haplotype information fits breeding values as random effects based on covariances defined by a “genomic relationship matrix” computed from genotypes [2]. These two models have been shown to be equivalent in terms of predicting breeding values [3, 4]

and we refer to them here as marker effect models (MEM) or breeding value models (BVM), the latter often known as Genomic Best Linear Unbiased Prediction (GBLUP).

Cross validation is often used to quantify the predictive ability of a statistical model. In  $k$ -fold cross validation, the whole dataset is partitioned into  $k$  parts with  $k$  analyses, where one part is omitted for training with validation on the omitted part. Leave-one-out cross validation (LOOCV) is a special case of  $k$ -fold cross validation with  $k = n$ , the number of observations. When the dataset is small, leave-one-out cross validation is appealing as the size of the training set is maximized. However, naive application of LOOCV is computationally intensive, requiring  $n$  analyses.

We show below how LOOCV can be performed using either the MEM or BVM with little more effort than is required for a single analysis with  $n$  observations.

## Methods

Use of the MEM is more efficient when the number  $n$  of individuals is larger than the number  $p$  of markers,

\*Correspondence: rohan@iastate.edu

<sup>1</sup>Department of Animal Science, Iowa State University, 50011, Ames, Iowa, USA  
Full list of author information is available at the end of the article

because for this model the mixed model equations are of order  $p$  plus the number of other effects. When  $n < p$ , estimated breeding values can be obtained more efficiently by solving the mixed model equations for the BVM of order  $n$  plus the number of other effects. We deal with the special case where the only other effect is a general mean and phenotypes have been pre-corrected for other nuisance variables. Efficient strategies for LOOCV using this special case for MEM when  $n \geq p$  and BVM when  $p \geq n$  are shown below.

**Marker effect models**

The MEM for GBLUP can be written as

$$y = \mathbf{1}\mu + X\beta + e, \tag{1}$$

where  $y$ , a  $n \times 1$  vector for phenotypes, has been pre-corrected for all fixed effects other than  $\mu$ , the overall mean,  $X$  is the  $n \times p$  matrix of marker covariates,  $\beta$  is a  $p \times 1$  random vector of the allele substitution effects and  $e$  is a  $n \times 1$  random vector of residuals. Often it is assumed that marker effects are identically and independently distributed (iid) random variables with null means and variances  $\sigma_\beta^2$ . Thus, under the usual assumption that the residuals are iid with null means and variances  $\sigma_e^2$ ,  $E(y) = \mathbf{1}\mu$ . When MEM is used, LOOCV can be performed by using a well-known strategy used in least-squares regression to compute the predicted residual sum of square (PRESS) [5] statistic.

**LOOCV strategy for MEM**

BLUP of  $\beta^* = \begin{bmatrix} \mu \\ \beta \end{bmatrix}$  can be obtained by solving the mixed model equations

$$(X^{*'}X^* + D\lambda)\hat{\beta}^* = X^{*'}y, \tag{2}$$

where  $X^* = [\mathbf{1} \ X]$ ,  $\hat{\beta}^* = \begin{bmatrix} \hat{\mu} \\ \hat{\beta} \end{bmatrix}$ ,  $D$  is a diagonal matrix whose elements are 0 followed by a  $p$  vector of 1s and  $\lambda = \frac{\sigma_e^2}{\sigma_\beta^2}$ .

Now, BLUP for  $\beta_{-j}^*$ , where observation  $j$  is left out, can be obtained as

$$\hat{\beta}_{-j}^* = (X_{-j}^{*'}X_{-j}^* + D\lambda)^{-1} X_{-j}^{*'}y_{-j}, \tag{3}$$

where  $X_{-j}^*$  is  $X^*$  with the  $j$ th row removed and  $y_{-j}$  is  $y$  with the  $j$ th element removed.

Suppose  $x_j^{*'}$  is the  $j$ th row of  $X^*$ , then from the matrix inverse lemma [4],

$$\begin{aligned} (X_{-j}^{*'}X_{-j}^* + D\lambda)^{-1} &= (X^{*'}X^* + D\lambda - x_j^{*'}x_j^*)^{-1} \\ &= (X^{*'}X^* + D\lambda)^{-1} \\ &\quad - \frac{(X^{*'}X^* + D\lambda)^{-1}x_j^{*'}x_j^*(X^{*'}X^* + D\lambda)^{-1}}{1 - H_{jj}}, \end{aligned} \tag{4}$$

where the quadratic  $H_{jj} = x_j^{*'}(X^{*'}X^* + D\lambda)^{-1}x_j^*$  is the  $j$ th diagonal element of  $H = X^*(X^{*'}X^* + D\lambda)^{-1}X^{*'}$ .

Using (3) in (4), the prediction residual for the  $j$ th observation can be written as

$$\begin{aligned} \hat{e}_j &= y_j - x_j^{*'}\hat{\beta}_{-j}^* \\ &= y_j - x_j^{*'} \left[ (X^{*'}X^* + D\lambda)^{-1} \right. \\ &\quad \left. - \frac{(X^{*'}X^* + D\lambda)^{-1}x_j^{*'}x_j^*(X^{*'}X^* + D\lambda)^{-1}}{1 - H_{jj}} \right] X_{-j}^{*'}y_{-j} \\ &= \frac{(1 - H_{jj})y_j - x_j^{*'}(X^{*'}X^* + D\lambda)^{-1}X_{-j}^{*'}y_{-j}}{1 - H_{jj}} \end{aligned} \tag{5}$$

$$= \frac{y_j - x_j^{*'}\hat{\beta}^*}{1 - H_{jj}}. \tag{6}$$

These prediction errors can be squared and accumulated over  $n$  realizations to compute PRESS defined as  $\sum_{j=1}^n \hat{e}_j^2$ . The accuracy of genomic prediction is often quantified as the correlation between the predicted and observed values of  $y_j$ , and that correlation can be estimated from the values of  $\hat{y}_j$ , which can be computed efficiently as  $\hat{y}_j = y_j - \hat{e}_j$ , using the observed values of  $y_j$ . When a specific group of individuals is of interest, prediction accuracies and PRESS can also be calculate using  $\hat{e}_j$  for individuals in that group.

**Breeding value models**

When  $n < p$ , the genomic prediction of the breeding value  $x_j'\hat{\beta}$  can be obtained more efficiently by solving the mixed model equations for the BVM:

$$y = \mathbf{1}\mu + Z\mathbf{u} + e, \tag{7}$$

where  $\mathbf{u} = X\beta$ ,  $\text{var}(\mathbf{u}) = \mathbf{X}\mathbf{X}'\sigma_\beta^2$ ,  $Z$  is the identity matrix of order  $n$  and other variables are as in the MEM. Further, in both models  $E(y) = \mathbf{1}\mu$ , and  $\text{var}(y) = \mathbf{X}\mathbf{X}'\sigma_\beta^2 + \mathbf{I}\sigma_e^2$ . These two models are said to be equivalent [6], and linear functions predicted from one model are identical to corresponding predictions from the other model. Two efficient strategies for LOOCV using the BVM are shown below.

**LOOCV strategy I for BVM**

The mixed model equations for this model are:

$$(\mathbf{Z}^{*'}\mathbf{Z}^* + \mathbf{G}\lambda)\hat{\mathbf{u}}^* = \mathbf{Z}^{*'}\mathbf{y}, \tag{8}$$

where  $\mathbf{Z}^* = [\mathbf{1} \ \mathbf{Z}]$ ,  $\hat{\mathbf{u}}^* = \begin{bmatrix} \hat{\mu} \\ \hat{\mathbf{u}} \end{bmatrix}$ ,  $\mathbf{G} = \begin{bmatrix} 0 & \mathbf{0}' \\ \mathbf{0} & (\mathbf{X}\mathbf{X}')^{-1} \end{bmatrix}$ .

Due to the relative order of the coefficient matrices for the MEM and the BVM, when  $n < p$ ,  $\mathbf{x}_j^{*'}\hat{\boldsymbol{\beta}}^*$  is more efficiently obtained as  $\hat{u}_j^*$ . Similarly,  $\text{var}(\mathbf{x}_j^{*'}\hat{\boldsymbol{\beta}}^* - \mathbf{x}_j^{*'}\hat{\boldsymbol{\beta}}^*) = \mathbf{x}_j^{*'}(\mathbf{X}^{*'}\mathbf{X}^* + \mathbf{D}\lambda)^{-1}\mathbf{x}_j^*\sigma_e^2$  can be obtained more efficiently as  $\text{var}(u_j^* - \hat{u}_j^*) = \mathbf{z}_j^{*'}(\mathbf{Z}^{*'}\mathbf{Z}^* + \mathbf{G}\lambda)^{-1}\mathbf{z}_j^*\sigma_e^2$ . Using these two equalities, the formula for  $\hat{e}_j$  becomes:

$$\begin{aligned} \hat{e}_j &= y_j - \mathbf{x}_j^{*'}\hat{\boldsymbol{\beta}}^*_{-j} \\ &= \frac{y_j - \mathbf{x}_j^{*'}\hat{\boldsymbol{\beta}}^*}{1 - H_{jj}} \\ &= \frac{y_j - \mathbf{x}_j^{*'}\hat{\boldsymbol{\beta}}^*}{1 - \mathbf{x}_j^{*'}(\mathbf{X}^{*'}\mathbf{X}^* + \mathbf{D}\lambda)^{-1}\mathbf{x}_j^*} \\ &= \frac{y_j - \mathbf{z}_j^{*'}\hat{\mathbf{u}}^*}{1 - \mathbf{z}_j^{*'}(\mathbf{Z}^{*'}\mathbf{Z}^* + \mathbf{G}\lambda)^{-1}\mathbf{z}_j^*} \\ &= \frac{y_j - \mathbf{z}_j^{*'}\hat{\mathbf{u}}^*}{1 - C_{jj}}, \end{aligned} \tag{9}$$

where the quadratic  $\mathbf{z}_j^{*'}(\mathbf{Z}^{*'}\mathbf{Z}^* + \mathbf{G}\lambda)^{-1}\mathbf{z}_j^*$  is the  $j$ th diagonal element of  $\mathbf{C} = \mathbf{Z}^*(\mathbf{Z}^{*'}\mathbf{Z}^* + \mathbf{G}\lambda)^{-1}\mathbf{Z}^{*}$ .

**LOOCV strategy II for BVM**

Another efficient strategy for BVM is shown here. First we consider the situation where  $\mathbf{y}$  has been pre-corrected for  $\mu$  in addition to nuisance effects so that  $E(\mathbf{y}) = \mathbf{0}$  and we define  $\text{var}(\mathbf{y}) = \mathbf{X}\mathbf{X}'\sigma_\beta^2 + \mathbf{I}\sigma_e^2 = \mathbf{V}$ . Now matrix  $\mathbf{Q}$  is constructed by augmenting the covariance matrix of  $\mathbf{y}$  with one leading row and column as

$$\mathbf{Q} = \begin{bmatrix} \mathbf{y}'\mathbf{y} & \mathbf{y}' \\ \mathbf{y} & \mathbf{V} \end{bmatrix}.$$

To obtain the prediction error for observation  $j$ , the second row and column of  $\mathbf{Q}$  are permuted with row and column  $j + 1$ . In this manner  $\mathbf{Q}$  has its rows and columns symmetrically permuted as  $\mathbf{P}'_j\mathbf{Q}\mathbf{P}_j = \mathbf{W}$ , where the permutation matrix  $\mathbf{P}_j$  is obtained by permuting the second row of the  $n$  order identity matrix with row  $j + 1$ . So the permuted matrix is:

$$\mathbf{W} = \begin{bmatrix} \mathbf{y}'\mathbf{y} & y_j & \mathbf{y}'_{-j} \\ y_j & V_{jj} & \mathbf{V}_{j,-j} \\ \mathbf{y}'_{-j} & \mathbf{V}_{-j,j} & \mathbf{V}_{-j,-j} \end{bmatrix} = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}' & \mathbf{C} \end{bmatrix}$$

where we will define the leading  $2 \times 2$  matrix as  $\mathbf{A} = \begin{bmatrix} \mathbf{y}'\mathbf{y} & y_j \\ y_j & V_{jj} \end{bmatrix}$ , and the other partitions as  $\mathbf{B} = \begin{bmatrix} \mathbf{y}'_{-j} \\ \mathbf{V}_{j,-j} \end{bmatrix}$ , and  $\mathbf{C} = \mathbf{V}_{-j,-j}$ , where  $-j$  denotes that the  $j$ th element, row or column has been removed. Defining  $\mathbf{W}^{11}$  as the top left or leading  $2 \times 2$  sub-matrix in  $\mathbf{W}^{-1}$  corresponding to the position of  $\mathbf{A}$  in  $\mathbf{W}$ , and using partitioned inverse-matrix identities [7], the inverse of  $\mathbf{W}^{11}$  can be written as,

$$\begin{aligned} (\mathbf{W}^{11})^{-1} &= \mathbf{A} - \mathbf{B}\mathbf{C}^{-1}\mathbf{B}' \\ &= \begin{bmatrix} \mathbf{y}'\mathbf{y} & y_j \\ y_j & V_{jj} \end{bmatrix} - \begin{bmatrix} \mathbf{y}'_{-j} \\ \mathbf{V}_{j,-j} \end{bmatrix} \mathbf{V}_{-j,-j}^{-1} \begin{bmatrix} \mathbf{y}_{-j} & \mathbf{V}_{-j,j} \end{bmatrix} \\ &= \begin{bmatrix} \mathbf{y}'\mathbf{y} - \mathbf{y}'_{-j}\mathbf{V}_{-j,-j}^{-1}\mathbf{y}_{-j} & y_j - \mathbf{y}'_{-j}\mathbf{V}_{-j,-j}^{-1}\mathbf{V}_{-j,j} \\ y_j - \mathbf{V}_{j,-j}\mathbf{V}_{-j,-j}^{-1}\mathbf{V}_{-j,j} & V_{jj} - \mathbf{V}_{j,-j}\mathbf{V}_{-j,-j}^{-1}\mathbf{V}_{-j,j} \end{bmatrix}. \end{aligned} \tag{10}$$

Now  $\mathbf{V}_{j,-j}$  in element (2, 1) of the above inverse matrix is the vector of covariances between  $y_j$  and  $\mathbf{y}_{-j}$  and  $\mathbf{V}_{-j,-j}^{-1}$  is the inverse of the covariance matrix of  $\mathbf{y}_{-j}$ . Thus,  $\hat{y}_j = \mathbf{V}_{j,-j}\mathbf{V}_{-j,-j}^{-1}\mathbf{y}_{-j}$  is the Best Linear Predictor (BLP) of  $y_j$  given  $\mathbf{y}_{-j}$ , and element (2,1) of (10) is the prediction error of  $y_j$ . The element (2,2) in (10) is the prediction error variance (PEV) for  $y_j$ , where  $PEV = \text{var}(y_j - \hat{y}_j)$ . PEV can also be used to calculate theoretical reliability for individual  $i$  as  $1 - \frac{PEV_i}{V_{jj}}$ , and characterizing the distributions of reliability for all the individuals in a dataset has a number of practical applications. Note this allows us to obtain the PEV of every individual and the distribution of these values provide information as to the robustness of genomic predictions across the population of individuals represented in the dataset. This PEV is determined by the genomic variance-covariance matrix and does not depend on  $\mathbf{y}$ . Two different datasets could generate the same PRESS statistic but with different distributions of PEV.

Now, because the permutation matrix  $\mathbf{P}_j$  is orthogonal,  $\mathbf{W}^{-1} = (\mathbf{P}'_j\mathbf{Q}\mathbf{P}_j)^{-1} = \mathbf{P}'_j\mathbf{Q}^{-1}\mathbf{P}_j$ , and the elements of  $\mathbf{W}^{11}$  that are of interest in terms of predicting individual  $j$  can be obtained directly from  $\mathbf{Q}^{-1}$  as

$$\mathbf{W}^{11} = \begin{bmatrix} q^{1,1} & q^{1,(1+j)} \\ q^{(1+j),1} & q^{(1+j),(1+j)} \end{bmatrix}. \tag{11}$$

It follows that  $\hat{e}_j$ , which is the off-diagonal element of the inverse of the  $2 \times 2$  matrix  $\mathbf{W}^{11}$ , can be written in terms of  $\mathbf{Q}^{-1}$  as

$$\hat{e}_j = \frac{-q^{(1+j),1}}{q^{1,1}q^{(1+j),(1+j)} - q^{1,(1+j)}q^{(1+j),1}}, \tag{12}$$

where  $q^{ij}$  is the element from row  $i$  and column  $j$  of  $\mathbf{Q}^{-1}$ . Thus, once  $\mathbf{Q}^{-1}$  is computed,  $\hat{e}_j$  for all  $j$  can be com-

**Table 1** Phenotypes and genotypes at 5 markers for 3 individuals used in the numerical example

	M1	M2	M3	M4	M5	Phenotypes
1	1	2	1	2	2	1.97
2	2	1	0	1	1	2.12
3	0	0	2	1	2	-0.62

puted using (12), and these values can be used to compute PRESS as  $\sum_{j=1}^n \hat{e}_j^2$ . To estimate the correlation between the predicted and observed values of  $y_j$ , the value of  $\hat{y}_j$  is efficiently obtained as the difference  $\hat{y}_j = y_j - \hat{e}_j$ .

Now we consider the situation without pre-correcting  $y$  for  $\mu$ , where  $E(\mathbf{y}) = \mathbf{1}\mu$ . Now the mixed model (7) contains both fixed and random effects. Note that the mixed model equations that correspond to this mixed effects model can be derived by treating  $\mu$  as “random” with null mean and large variance. So, let

$$var(\boldsymbol{\beta}^*) = \begin{bmatrix} \sigma_L^2 & \mathbf{0}' \\ \mathbf{0} & I\sigma_\beta^2 \end{bmatrix} = \boldsymbol{\Sigma},$$

for sufficiently large value of  $\sigma_L^2$ . Then under this assumption,  $E(\mathbf{y}) = \mathbf{0}$  and  $var(\mathbf{y}) = \mathbf{X}^*\boldsymbol{\Sigma}\mathbf{X}^{*'} + \mathbf{I}\sigma_e^2 = \mathbf{V}^*$ , and thus  $\hat{y}_j = \mathbf{V}_{j,-j}^* \mathbf{V}_{-j,-j}^{*-1} \mathbf{y}_{-j}$  is the BLP from the random effects rather than mixed effects model of  $y_j$  given  $\mathbf{y}_{-j}$ . This BLP obtained from the model with random  $\mu$  will be numerically very close to the BLUP obtained from the mixed model with fixed  $\mu$ . The  $\mathbf{Q}$  matrix corresponding to the BLP with random  $\mu$  is constructed as  $\begin{bmatrix} \mathbf{y}'\mathbf{y} & \mathbf{y}' \\ \mathbf{y} & \mathbf{V}^* \end{bmatrix}$  and prediction residuals are obtained as (12).

**Numerical example**

Phenotypes  $y$  and genotypes  $X$  at 5 markers for 3 individuals are in Table 1. Assume  $\sigma_\beta^2 = \frac{\sigma_e^2}{10}$  and the overall mean  $\mu$  is the only fixed effect. In LOOCV strategy for MEM and strategy I for BVM, the diagonal elements of  $\mathbf{H}$  for MEM and  $\mathbf{C}$  for BVM, which are in the denominators of (6) and (9), are in Table 2. The numerators of (6) and (9) are obtained by solving the MME (2) and (8). Then prediction errors are calculated as in (6) and (9) and shown in Table 4. In LOOCV strategy II for BVM, the  $\mathbf{Q}$  matrix (Table 3) is constructed using  $\sigma_L^2 = 1000$ , which is sufficiently large

**Table 2** Diagonal elements of  $\mathbf{H}$  in LOOCV strategy for MEM and  $\mathbf{C}$  for BVM

	$j = 1$	$j = 2$	$j = 3$
$H_{jj}$	0.46	0.51	0.55
$C_{jj}$	0.46	0.51	0.55

**Table 3**  $\mathbf{Q}$  matrix in strategy II for BVM

	1	2	3	4
1	8.75	1.97	2.12	-0.62
2	1.97	1,002.40	1,000.80	1,000.80
3	2.12	1,000.80	1,001.70	1,000.30
4	-0.62	1,000.80	1,000.30	1,001.90

relative to  $\sigma_e^2$  for  $\mu$  to be indistinguishable from a fixed effect with a flat prior. The prediction errors are calculated as (12) and shown in Table 4. The MEM strategy and BVM strategy I gave identical prediction errors and identical PRESS for this numerical example were numerically very close to those from the BVM strategy II.

**Simulation to compare efficiency**

Two datasets were simulated using XSim [8], where 1,000 offspring were sampled from random mating of 100 parents for 10 non-overlapping generations, to compare the computational efficiencies for naive and efficient strategies using BVM or MEM for LOOCV in GBLUP. Dataset I was simulated with 1,000 observations and 10,000 SNP markers for a  $p \gg n$  scenario. Dataset II was simulated with 1,000 observations and 100 markers for a  $n \gg p$  scenario. The processor used in the analyses was a 1.4 GHz Intel Core i5 with 4 GB of memory.

For dataset II, efficient MEM is 99 times faster than the naive application (2.979 s versus 0.030 s) (Table 5). All strategies implemented in Julia, a scientific programming language, gave virtually identical prediction accuracies defined as the correlation between  $\mathbf{y}$  and  $\hat{\mathbf{y}}$  for each dataset. For dataset I, efficient BVM is 786 times faster than the naive application (3.107 s versus 2,442.59 s) (Table 5).

**Discussion**

In genomic prediction, the candidates to be predicted are often offspring that are genotyped but not yet phenotyped. In this situation, LOOCV using all individuals in the training dataset will provide an upper bound for the accuracy of prediction, because ancestors in the training dataset with large numbers of descendants have more accurate predictions than descendants. A better estimate of the accuracy of prediction can be obtained by applying LOOCV to only terminal offspring in the training dataset.

**Table 4** Prediction errors from different LOOCV strategies (different strategies gave identical prediction errors)

	$j = 1$	$j = 2$	$j = 3$
$\hat{e}_j$	1.13	1.21	-2.66

**Table 5** Efficiency of alternative LOOCV strategies for GBLUP

	Alternative LOOCV strategies				
	Naive MEM	Naive BVM	Efficient MEM	Efficient BVM I	Efficient BVM II
$n = 1,000; p = 10,000$	9,490.608	2,442.590	105.141	3.107	5.945
$n = 1,000; p = 100$	2.979	169.928	0.030	2.725	0.217

Results are given for the computing time in seconds using naive MEM, naive BVM, efficient MEM, efficient BVM I and efficient BVM II

## Conclusions

Efficient strategies for LOOCV in GBLUP are presented in this paper. LOOCV strategy I and II for BVM are more efficient when  $p \gg n$ . LOOCV strategy for MEM is more efficient when  $n \gg p$ . The accuracy of genomic prediction is often quantified as the correlation between the predicted and observed values of  $y_j$ , and this correlation can be estimated efficiently using LOOCV strategies. Compared to naive application of LOOCV, which is computationally intensive, LOOCV can be implemented efficiently.

## Acknowledgements

Not applicable.

## Funding

This work was supported by the US Department of Agriculture, Agriculture and Food Research Initiative National Institute of Food and Agriculture Competitive grant no. 2015-67015-22947. The funders had no role in the design of the study and collection, analysis, and interpretation of data and in writing of the manuscript.

## Availability of data and materials

The datasets during and/or analysed during the current study available from the corresponding author on reasonable request.

## Authors' contributions

All authors contributed to the development of the statistical methods. HC wrote the program code and conducted the analyses. The manuscript was prepared by HC, RLF and DG. All authors read and approved the final manuscript.

## Competing interests

The authors declare that they have no competing interests.

## Consent for publication

Not applicable.

## Ethics approval and consent to participate

Not applicable.

## Author details

<sup>1</sup>Department of Animal Science, Iowa State University, 50011, Ames, Iowa, USA. <sup>2</sup>Department of Statistics, Iowa State University, 50011, Ames, Iowa, USA. <sup>3</sup>Institute of Veterinary, Animal & Biomedical Science, Massey University, Palmerston North, New Zealand.

Received: 21 September 2016 Accepted: 27 March 2017

Published online: 02 May 2017

## References

1. Meuwissen THE, Hayes BJ, Goddard ME. Prediction of total genetic value using genome-wide dense marker maps. *Genetics*. 2001;157:1819–29.
2. Nejati-Javaremi A, Smith C, Gibson JP. Effect of total allelic relationship on accuracy of evaluation and response to selection. *J Anim Sci*. 1997;75:1738–45.

3. Fernando RL. Genetic evaluation and selection using genotypic, phenotypic and pedigree information. *Proc 6th Wld Cong Genet Appl Livest Prod*. 1998;26:329–36.
4. Strandén I, Garrick DJ. Technical note: Derivation of equivalent computing algorithms for genomic predictions and reliabilities of animal merit. *J Dairy Sci*. 2009;92(6):2971–75.
5. Allen DM. The relationship between variable selection and data augmentation and a method for prediction. *Technometrics*. 1974;16(1):125–7.
6. Henderson CR. *Applications of Linear Models in Animal Breeding*. Guelph: Univ. Guelph; 1984.
7. Searle SR. *Matrix Algebra Useful for Statistics*. New York: John Wiley and Sons, Inc; 1982.
8. Cheng H, Garrick D, Fernando R. Xsim: Simulation of descendants from ancestors with sequence data. *G3*. 2015;5(7):1415–17.

Submit your next manuscript to BioMed Central and we will help you at every step:

- We accept pre-submission inquiries
- Our selector tool helps you to find the most relevant journal
- We provide round the clock customer support
- Convenient online submission
- Thorough peer review
- Inclusion in PubMed and all major indexing services
- Maximum visibility for your research

Submit your manuscript at  
www.biomedcentral.com/submit

