

Using Machine Learning To Improve the Accuracy of Genomic Prediction on Reproduction Traits in Pigs

Xue Wang

China Agricultural University

Shaolei Shi

China Agricultural University

Guijiang Wang

Hebei Province Animal Husbandry and Improved Breeds Work station

Wenxue Luo

Hebei Province Animal Husbandry and Improved Breeds Work Station

Xia Wei

Zhangjiakou Dahao Heshan New Agricultural Development Co., Ltd

Ao Qiu

China Agricultural University

Fei Luo

Hebei Province Animal Husbandry and Improved Breeds Work Station

Xiangdong Ding (✉ xding@cau.edu.cn)

China Agricultural University

Research

Keywords: machine learning, genomic prediction, prediction accuracy, pig

Posted Date: November 19th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-1083849/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Using machine learning to improve the accuracy of genomic**
2 **prediction on reproduction traits in pigs**

3 Xue Wang^a, Shaolei Shi^a, Guijiang Wang^b, Wenxue Luo^b, Xia Wei^c, Ao Qiu^a, Fei
4 Luo^b, Xiangdong Ding^{a*}

5 a Key Laboratory of Animal Genetics and Breeding of Ministry of Agriculture and
6 Rural Affairs, National Engineering Laboratory of Animal Breeding, College of Animal
7 Science and Technology, China Agricultural University, Beijing, China

8 b Hebei Province Animal Husbandry and Improved Breeds Work Station, Shijiazhuang,
9 Hebei, China

10 c Zhangjiakou Dahao Heshan New Agricultural Development Co., Ltd, Zhangjiakou,
11 Hebei, China

12 * Corresponding author.

13 E-mail addresses:

14 xwangchnm@163.com

15 BS20183040372@cau.edu.cn

16 13933077085@126.com

17 13832350632@163.com

18 weixia16888@163.com

19 956482319@qq.com

20 43068883@qq.com

21 xding@cau.edu.cn

22 **Abstract**

23 **Background:** Recently, machine learning (ML) is becoming attractive in genomic
24 prediction, while its superiority in genomic prediction and the choosing of optimal ML
25 methods are needed investigation.

26 **Results:** In this study, 2566 Chinese Yorkshire pigs with reproduction traits records
27 were used, they were genotyped with GenoBaits Porcine SNP 50K and PorcineSNP50
28 panel. Four ML methods, including support vector regression (SVR), kernel ridge
29 regression (KRR), random forest (RF) and Adaboost.R2 were implemented. Through
30 20 replicates of five-fold cross-validation, the genomic prediction abilities of ML
31 methods were explored. Compared with genomic BLUP(GBLUP), single-step GBLUP
32 (ssGBLUP) and Bayesian method BayesHE, our results indicated that ML methods
33 significantly outperformed. The prediction accuracy of ML methods was improved by
34 19.3%, 15.0% and 20.8% on average over GBLUP, ssGBLUP and BayesHE, ranging
35 from 8.9% to 24.0%, 7.6% to 17.5% and 11.1% to 24.6%, respectively. In addition, ML
36 methods yielded smaller mean squared error (MSE) and mean absolute error (MAE) in
37 all scenarios. ssGBLUP yielded improvement of 3.7% on average compared to GBLUP,
38 and the performance of BayesHE was close to GBLUP. Among four ML methods, SVR
39 and KRR had the most robust prediction abilities, which yielded higher accuracies,
40 lower bias, lower MSE and MAE, and comparable computing efficiency as GBLUP.
41 RF demonstrated the lowest prediction ability and computational efficiency among ML
42 methods.

43 **Conclusion:** Our findings demonstrated that ML methods are more efficient than

44 traditional genomic selection methods, and it could be new options for genomic
45 prediction.

46 **Key words:** machine learning, genomic prediction, prediction accuracy, pig

47

48 **Background**

49 Genomic selection (GS) has been widely recognized and successfully implemented in
50 animal and plant breeding programs ^[1-3]. It is reported that the breeding costs of dairy
51 cattle using GS were 92% lower than that of tradition progeny testing ^[4]. At present, the
52 genetic gain rate of the annual yield traits of US Holstein dairy cattle has increased from
53 around 50% to 100%^[5]. The accuracy of GS depends on methods of genomic breeding
54 values estimation (GEBV), reference population size, marker density, and heritability,
55 *etc.* Currently, parametric methods are most commonly used for livestock and poultry
56 genomic selection, mainly based on either the genomic covariance between genotyped
57 individuals *e.g.* genomic BLUP (GBLUP)^[6] or single-step GBLUP (ssGBLUP)^[7, 8]) or
58 Bayesian regression models^[9, 10], with differences mainly depends on the prior
59 distribution of marker effects. Nevertheless, these linear models usually only take into
60 account the additive inheritance and ignore the complex non-linear relationships that
61 may exist between markers and phenotypes (*e.g.* epistasis, dominance, genotype-by-
62 environment interactions). In addition, parametric methods usually provide limited
63 flexibility for handling non-linear effects in high-dimensional genomic data, resulting
64 in huge computational demands ^[11], while considering nonlinearity may enhance the
65 predictive ability of complex traits ^[12]. Therefore, new strategy should be explored to

66 more accurately estimate the genomic breeding values.

67 Driven by applications in intelligent robots, self-driving cars, automatic translation,
68 face recognition, artificial intelligence games and medical services, machine learning
69 (ML) has gained considerable attention in the past decade. Some characteristics of the
70 ML methods make it potentially attractive to deal with high-order non-linear
71 relationships in high-dimensional genomic data, *e.g.* allowing the number of variables
72 larger than the sample size ^[13], capable of capturing the hidden relationship between
73 genotype and phenotype in an adaptive manner, and imposing little or no specific
74 distribution assumptions about the predictor variables as GBLUP and Bayesian
75 methods ^[14, 15] .

76 Studies have shown that random forest (RF), support vector regression (SVR), kernel
77 ridge regression (KRR) and other machine learning methods gained advantage over
78 GBLUP and Bayes B, *etc.* ^[16-18]. Ornella et al. compared the performance of support
79 vector regression, random forest regression, Reproducing Kernel Hilbert space (RKHS),
80 ridge regression, and Bayesian Lasso in genomic prediction, and found that RKHS and
81 random forest regression were the best ^[19]. González-Camacho et al. reported the
82 support vector machine (SVM) with linear kernel performed the best in comparison
83 with other ML methods and linear models in the genomic prediction of the rust
84 resistance of wheat ^[18]. Additionally, ML algorithms have also been widely used in the
85 fields of gene screening, genotype imputation, and protein structure and function
86 prediction, *etc.* ^[20-23], demonstrating its superiority as well. However, one challenge for
87 the ML is choosing the optimum ML method as a series of ML algorithms have been

88 proposed and each has its own characteristics and shows different prediction abilities
89 in different datasets and traits.

90 Therefore, the objective of this study was to assess the performance of machine learning
91 methods in genomic prediction through the comparison with existing prevail GBLUP
92 and Bayesian methods, and on the other hand, the efficiency of different ML methods
93 were compared as well in order to explore the ideal ML algorithm for genomic
94 prediction.

95 **Materials and Methods**

96 *Ethics Statement*

97 The whole procedure for blood sample collection was carried out in strict accordance
98 with the protocol approved by the Animal Care and Use Committee of China
99 Agricultural University (Permit Number: DK996).

100 *Population and Phenotypes*

101 A Yorkshire pig population from DHHS, a breeding farm in Hebei province, China, was
102 studied. A total 2566 animals born between 2016 and 2020 were sampled and 4274
103 reproductive records of total number of piglets born (TNB) and number of piglets born
104 alive (NBA) were available, and 3893 animals were traced back to construct pedigree
105 relationship matrix (A matrix). A single-trait repeatability model was used to estimate
106 heritabilities. The fixed effects included herd-year-season, and the random effects
107 included additive genetic effects, random residuals, and permanent effects. The
108 information of the animals, phenotypes and genetic components, as well as the
109 estimated heritabilities were listed in Table 1. The estimated heritabilities of TNB and

110 NBA were both 0.12.

111 *Derivation of corrected phenotypes*

112 In order to avoid double counting of parental information, the corrected phenotypes (y_c)
113 derived from the estimated breeding values (EBV) were used as response variable in
114 genomic prediction. The pedigree-based BLUP and single-trait repeatability model
115 were performed to estimate the breeding values for each trait separately.

$$116 \quad y = Xb + Z_a a + Z_{pe} pe + e, \quad (1)$$

117 where y is the vector of phenotypic values; b is the vector of fixed effects including
118 herd-year-season; a represent additive genetic effects, following a norm distribution
119 $N(0, A\sigma_a^2)$, where A is the pedigree-based relationship matrix, σ_a^2 is the additive
120 genetic variance. pe is permanent environment effects with norm distribution $N(0,$
121 $I\sigma_{pe}^2)$, where σ_{pe}^2 is permanent environment variance. e is the vector of random error,
122 following a norm $N(0, I\sigma_e^2)$, where σ_e^2 represents residual variance. X , Z_a , and Z_{pe}
123 are incidence matrices linked b , a and pe to y . A total of 3893 individuals were
124 traced to construct A matrix. Their EBVs were calculated using the DMUAI procedure
125 of the DMU software [24]. The y_c were calculated as EBV plus the average estimated
126 residuals for multiple parties of a sow following Guo et al. [25].

127 *Genotype data and imputation*

128 Two kinds of 50K SNP panels, PorcineSNP50 BeadChip (Illumina, CA, USA) and
129 GenoBaits Porcine SNP 50K (Molbreeding, China) were used for the genotyping. A
130 total of 1189 sows were genotyped with PorcineSNP50 BeadChip, which includes
131 50,697 SNPs across the genome, and 1978 individuals were genotyped using GenoBaits

132 Porcine SNP 50K with 52,000 SNPs. There are 30,998 common SNPs between these
133 two SNP panels, and 601 individuals were genotyped with both SNP panels and, 2566
134 genotyped individuals were therefore finally used for further analysis including 1189
135 animals with PorcineSNP50 BeadChip and 1377 pigs with GenoBaits Porcine SNP 50K.
136 The animals genotyped with GenoBaits Porcine SNP 50K were imputed to
137 PorcineSNP50 BeadChip using Beagle 5.0 [26]. The reference population size for
138 genotype imputation was 3720. Imputation accuracy was assessed by the dosage R-
139 squared measure (DR2), which is the estimated squared correlation between the
140 estimated allele dose and the true allele dose. The genotype correlation (COR) and the
141 genotype concordance rate (CR) were also calculated based on the 601 overlap animals
142 to evaluate the imputation accuracy. After imputation, the quality control on genotype
143 were carried out using PLINK software [27], SNPs with minor allele frequency (MAF)
144 lower than 0.01 and call rate lower than 0.90 were removed, and individuals with call
145 rate lower than 0.90 were excluded. Finally, all animals and 44,922 SNPs on autosomes
146 were remained for further analysis.

147 *Statistical models*

148 GBLUP, ssGBLUP, Bayesian Horseshoe (BayesHE) and four ML regression methods,
149 support vector regression (SVR), Kernel ridge regression (KRR), Random forest (RF),
150 and Adaboost.R2 were used to predict GEBV. For ssGBLUP, in order to prevent the
151 problem that singular matrix cannot be inverted, $G_w = (I-w)G_a + wA_{22}$, and w was equal
152 to 0.05 [28]. BayesHE was developed by Shi. et al [29], it is based on Global-local priors
153 to increase the flexibility and adaptability of the Bayesian model. In this study, the first

154 form of BayesHE (BayesHE1) was used ^[29], and the Markov chain Monte Carlo
155 (MCMC) chain was run for 50,000 cycles, with the first 20,000 cycles being discarded
156 as burn-in and every 50 sample of the remaining 30,000 iterations were saved to infer
157 posterior statistics. In-house scripts written in Fortran 95 were used for BayesHE
158 analyses, and the DMUAI procedure implemented in DMU software was used for
159 GBLUP and ssGBLUP analyses. Meanwhile, the four ML regression methods are
160 introduced as follows.

161 *Support vector regression*

162 Support vector machine (SVM) was proposed by Vapnik ^[30] for binary classification.
163 SVR is the application of SVM in regression for dealing with quantitative responses,
164 which uses a linear or non-linear kernel function to map the input space (the marker
165 dataset) to a higher dimensional feature space, and performed modeling and prediction
166 on the feature space ^[31]. In other words, we can build a linear model in the feature space
167 to deal with regression problems. The model formulation of SVR can be expressed as:

168
$$f(x) = \beta_0 + h(x)^T \beta, \quad (2)$$

169 in which $h(x)^T \beta$ is the kernel function, β is the vector of weights, and β_0 is the bias.
170 Generally, the formalized SVR is given by minimizing the following restricted loss
171 function:

172
$$\min_{\beta_0, \beta} \frac{1}{2} \|\beta\|^2 + C \sum_{i=1}^n V(y_i - f(x_i)), \quad (3)$$

173 in which

174
$$V_\varepsilon(r) = \begin{cases} 0, & \text{if } |r| < \varepsilon \\ |r| - \varepsilon, & \text{otherwise} \end{cases} \quad (4)$$

175 $V_\varepsilon(r)$ is the ε -insensitive loss and C (“cost parameter”) is the regularization constant

176 that controls the trade-off between prediction error and model complexity. y is a
 177 quantitative response and $\|\cdot\|$ is the norm in the Hilbert space. After optimization, the
 178 final form of SVR can be written as:

$$179 \quad f(x) = \sum_{i=1}^m (\hat{a}_i - a_i) k(x, x_i), \quad (5)$$

180 in which $k(x_i, x_j) = \phi(x_i)^T \phi(x_j)$ is the kernel function. In this research, grid search
 181 was used to find the best kernel function and the best hyper-parameters of C and
 182 gamma. An internal five-fold cross validation (5-fold CV) strategy was performed to
 183 adjust the hyper-parameters when performing a grid search.

184 *Kernel ridge regression*

185 Kernel ridge regression (KRR) is a non-linear regression method, which can effectively
 186 discover the non-linear structure of the data^[32]. KRR uses a non-linear kernel function
 187 to map the data to a higher dimensional kernel space, and then builds a ridge regression
 188 model to make the data linearly separable in this kernel space. The linear function in
 189 the kernel space is selected according to the mean squared error loss of ridge
 190 regularization^[32]. The final KRR prediction model can be written as:

$$191 \quad y(x_i) = k'(K + \lambda I)^{-1} \hat{y}, \quad (6)$$

192 where λ is the regularization constant; K is the Gram matrix with entries $K_{ij} =$
 193 $K(x_i, x_j) = \phi(x_i) \cdot \phi(x_j)^T$, thus, for n training samples, the obtained kernel matrix is:

$$194 \quad K = \begin{bmatrix} K(x_1, x_1) & K(x_1, x_2) & \Lambda & K(x_1, x_n) \\ K(x_2, x_1) & K(x_2, x_2) & \Lambda & K(x_2, x_n) \\ \mathbf{M} & \mathbf{M} & \mathbf{M} & \mathbf{M} \\ K(x_n, x_1) & K(x_n, x_2) & \Lambda & K(x_n, x_n) \end{bmatrix}_{n \times n}. \quad (7)$$

195 I is the identity matrix; $k' = K(x_i, x_j)$ with $j = 1, 2, 3, \dots, n$, n is the number of
 196 training samples, and x_i is the test sample. In the expanded form,

$$k' = \begin{bmatrix} K(x_i, x_1) \\ K(x_i, x_2) \\ \vdots \\ K(x_i, x_n) \end{bmatrix}. \quad (8)$$

The grid search was used to find the most suitable kernel function and λ in this study, and an internal 5-fold CV strategy was used for tuning the hyper-parameters.

Random forest

Random forest (RF) is a machine learning method that uses voting or the average of multiple decision trees to determine the classification or predicted values of new instances^[33]. Random forest is essentially a collection of decision trees, and each decision tree is slightly different from other trees^[34]. Random forest reduces the risk of overfitting by averaging the prediction results of many decision trees^[18]. Random forest regression can be written in the following form:

$$y = \frac{1}{M} \sum_{m=1}^M t_m(\psi_m(y: X)), \quad (9)$$

in which y is the predicted value of random forest regression, $t_m(\psi_m(y: X))$ is an individual regression tree, and M is the number of decision trees in the forest. The prediction is obtained by passing down the predictor variables in the flowchart of each tree, and the corresponding estimated value at the terminal node is used as the predicted value. Finally, the predictions of each tree in RF are averaged to calculate the final prediction of unobserved data. The grid search was used to find the most suitable hyper-parameter M and the maximum depth of the tree, and the inner 5-fold CV was performed to tune the hyper-parameters.

Adaboost.R2

Adaboost.R2^[35] is an ad hoc modification of Adaboost.R and is an extension of

218 Adaboost.M2 to deal with regression problems, which repeatedly uses a regression tree
 219 as a weak learner followed by increasing the weights of incorrectly predicted samples
 220 and decreasing the weights of correctly predicted samples. It builds a “committee” by
 221 integrating multiple weak learners, making its prediction effect better than those of
 222 weak learners ^[36]. Adaboost.R2 regression model can be written as:

$$223 \quad y = \inf \left[y \in Y: \sum_{t: f_t(x) \leq y} \log \frac{1}{\varepsilon_t} \geq \frac{1}{2} \sum_t \log \frac{1}{\varepsilon_t} \right], \quad (10)$$

224 in which y is the predicted GEBV, $f_t(x)$ is predicted value of the t-th weak learner,
 225 and ε_t is the error rate of $f_t(x)$, $\varepsilon_t = \bar{L}_t / (1 - \bar{L}_t)$ and the average loss $\bar{L}_t =$
 226 $\sum_{i=1}^m L_t(i) D_t(i)$, in which $L_t(i)$ is the error between the actual observation value and
 227 the predicted value of the i-th predicted individual, and $D_t(i)$ is the weight distribution
 228 of $f_t(x)$. After $f_t(x)$ is trained, the weight distribution $D_t(i)$ will become $D_{t+1}(i)$,

$$229 \quad D_{t+1}(i) = \frac{D_t(i) \beta_t^{(1-L_t(i))}}{Z_t}, \quad (11)$$

230 in which Z_t is a normalization factor chosen such that $D_{t+1}(i)$ will be a distribution.

231 In current study, SVR and KRR were respectively used as weak learners of
 232 Adaboost.R2.

233 For these four ML methods, the vectors of genotypes (coded as 0, 1, 2) were the input
 234 independent variables and y_c were used as response variables, and Sklearn package for
 235 Python (V0.22) was used for genomic prediction.

236 Meanwhile, the optimal hyper-parameters for SVR, KRR, RF and Adaboost.R2
 237 according to the grid search were shown in Table S1.

238 *Accuracy of genomic prediction*

239 Five-fold cross validation was used to estimate the accuracies of genomic prediction,

240 in which 2566 individuals were randomly split into five groups with 513 individuals
241 each. For each cross validation, four of the five groups were defined as reference
242 population, and the left one was treated as the validation population. The genotyped
243 reference and validation sets in each replicate of 5-fold CV were same for all methods,
244 and it should be noted that non-genotyped individuals were added in the reference
245 population in ssGBLUP. For all methods, the accuracy of genomic prediction was
246 calculated as the Pearson correlation between the GEBVs and corrected phenotypes y_c
247 in validation population. In addition, the prediction unbiasedness was also calculated
248 as the regression of y_c on the GEBVs of validation population. The 5-fold CV scheme
249 was repeated 20 times, and the overall prediction accuracy and unbiasedness was the
250 average of 20 replicates. The Hotelling-Williams Test ^[37] was performed to compare
251 the prediction accuracy of different methods.

252 Meanwhile, prediction ability metrics *e.g.* mean squared error (MSE) and mean
253 absolute error (MAE) were also used to evaluate the performance of regression models
254 in the present study. MSE can take both prediction accuracy and bias into account ^[38],
255 and the smaller the value of MSE, the better the accuracy of the model to describe the
256 experimental data is. MAE can better reflect the actual situation of the predicted value
257 error. Their formulas can be written as follows.

258
$$MSE = \frac{1}{m} \sum_{i=1}^m (f_i - y_i)^2, \text{ and}$$

259
$$MAE = \frac{1}{m} \sum_{i=1}^m |f_i - y_i|, \quad (12)$$

260 where m represents the number of animals in each cross-validation test fold of the 5-
261 fold CV, f is the vector of predicted values (GEBVs) and y is the vector of observed

262 values (y_c). The final MSE and MAE were the average of 20 replicates.

263 **Results**

264 *Genotype imputation accuracy*

265 Figure 1 illustrates the accuracy of imputing GenoBaits Porcine SNP 50K to
266 PorcineSNP50 BeadChip across minor allele frequency (MAF) intervals and
267 chromosomes. DR2, CR and COR were not sensitive to MAF except that COR was
268 lower when the MAF was less than 0.05 and in the range of 0.45 to 0.5 (Figure 1a).
269 DR2, CR and COR on each chromosome were 0.978~0.988, 0.984~0.988 and
270 0.957~0.972, respectively, and no significant differences were observed in DR2, CR
271 and COR between chromosomes (Figure 1b). In the same scenarios, the values of COR
272 were smaller than those of DR2, CR. The averaged DR2, CR and COR across all
273 variants were 0.984, 0.985 and 0.964, respectively, indicating the imputation is enough
274 accurate to analysis two SNP panel together.

275 *Accuracy of genomic prediction*

276 *Comparison of ML methods with (ss)GBLUP and BayesHE*

277 Table 2 shows the prediction accuracies and unbiasedness of machine learning methods,
278 (ss)GBLUP and BayesHE on traits of TNB and NBA. The accuracies of ML methods
279 were significantly higher than those of (ss)GBLUP and BayesHE. The improvements
280 of ML methods over GBLUP, ssGBLUP and BayesHE were 19.3%, 15.0% and 20.8%
281 on average, ranging from 8.9% to 24.0%, 7.6% to 17.5% and 11.1% to 24.6%,
282 respectively. For trait TNB, compared with GBLUP, the average accuracy of all ML
283 methods in this study has been improved, support vector regression (SVR) gained

284 improvement of 19.0% as same as Kernel ridge regression (KRR), Adaboost.R2 based
285 on SVR and KRR obtained the improvement of 18.1% and 17.7%, respectively, while
286 random forest (RF) yielded the lowest improvement of 8.9% advantage over GBLUP.
287 The similar advantage of ML were also over ssGBLUP, the improvements of SVR,
288 KRR, RF, Adaboost.R2_SVR and Adaboost.R2_KRR were 17.5%, 17.5%, 7.6%, 16.7%
289 and 16.3%, respectively. ML methods gained the largest advantage over BayesHE, the
290 accuracy from SVR, KRR, RF, Adaboost.R2_SVR and Adaboost.R2_KRR were
291 respectively improved by 21.4%, 21.4%, 11.1%, 20.6% and 20.2% compared with
292 BayesHE. For trait NBA, although ML methods still performed better than GBLUP,
293 ssGBLUP and BayesHE, Adaboost.R2_KRR gained the largest improvement in all
294 comparisons, and KRR obtained the second largest improvement. SVR and
295 Adaboost.R2 based on SVR yielded same improvements on GBLUP, ssGBLUP and
296 BayesHE. RF still gained the lowest improvement compared with other ML methods.
297 Meanwhile, GBLUP, ssGBLUP and BayesHE had similar performance, and no
298 statistical differences of prediction accuracy were found among them. Nevertheless,
299 ssGBLUP produced average improvement of 3.7% compared with GBLUP (1.2% for
300 TNB; 6.3% for NBA), while less bias was observed by GBLUP in all scenarios.
301 BayesHE yielded similar accuracy with GBLUP (0.243 and 0.248 for TNB; 0.207 and
302 0.208 for NBA), but the unbiasedness of BayesHE was much closer to 1 (1.015 for
303 TNB; 1.009 for NBA).

304 On the other hand, mean squared error (MSE) and mean absolute error (MAE) were
305 also used to assess the performance of different methods. As shown in Table 3, ML

306 methods were generally superior to GBLUP, ssGBLUP and BayesHE in terms of MSE
307 and MAE. For two reproduction traits TNB and NBA, all ML methods yielded lower
308 MSE and MAE than GBLUP, ssGBLUP and BayesHE. The performance of GBLUP,
309 ssGBLUP and BayesHE was very close, and ssGBLUP produced a bit lower MSE (5.26
310 for TNB; 3.95 for NBA) and MAE (1.748 for TNB; 1.532 for NBA) among these three
311 methods, while they were still higher than those obtained from RF, which performed
312 the worst among four ML methods, and generated 5.212 and 3.901 of MSE and 1.747
313 and 1.527 of MAE on TNB and NBA, respectively. Among ML models, the
314 performance of SVR and KRR was the best, and they yielded the smallest MSE and
315 MAE in all scenarios.

316 *Comparison between ML methods*

317 Table 2 and 3 indicates that ML methods performed better than GBLUP, ssGBLUP and
318 BayesHE. They also show RF had the lowest accuracy even though no significant
319 differences were observed among the ML methods in this study. The accuracies of SVR,
320 KRR, Adaboost.R2_SVR and Adaboost.R2_KRR were improved by an average of
321 5.8%, 6.2%, 5.5% and 6.1% compared to RF, ranging from 8.1% to 9.3% for TNB and
322 from 2.4% to 4.0% for NBA, respectively. For TNB, SVR and KRR showed the highest
323 accuracies (0.295 for both), and Adaboost.R2_KRR yielded the highest accuracies on
324 NBA (0.258). In the meantime, in the comparison of unbiasedness, SVR produced the
325 lowest genomic prediction bias, and the regression coefficient was close to 1.0, while
326 Adaboost.R2 method with both base learner SVR and KRR produced larger bias. As a
327 trade-off between accuracy and unbiasedness, SVR and KRR had the most robust

328 prediction ability, which also confirmed by the results of MSE and MAE, in which SVR
329 and KRR had the smallest MSE and MAE in all scenarios.

330 It should be noted that the better performance of ML methods was acquired by tuning
331 hyper-parameters (Table S1). Compared with using the default hyper-parameters, the
332 accuracy was improved by 14.3% on average from the ML methods with optimal hyper-
333 parameters (Table S2), the accuracy of SVR, KRR, RF and Adaboost.R2 with optimal
334 hyper-parameters gained improvements by 15.7%, 11.7%, 9.8% and 15.0% respectively
335 on the genomic prediction accuracies for TNB, and for NBA, the improvements were
336 13.4%, 15.3%, 10.2% and 23.4%, respectively. As for unbiasedness, except for SVR on
337 TNB, the unbiasedness of all ML methods using the default parameters was lower than
338 the unbiasedness using the optimal parameters.

339 *Computing time*

340 The computing time of each method is demonstrated in Table 4. Among all methods,
341 KRR was the fastest algorithm, it took an average of 1.16 minutes in each iteration of
342 cross-validation to complete the analysis, requiring considerably less time than GBLUP
343 (2.07 minutes) and ssGBLUP (3.23 minutes). The computing efficiency of SVR (5.28
344 minutes) and Adaboost.R2_KRR (5.16 minutes) were comparable with KRR, GBLUP
345 and ssGBLUP. However, RF (53.45 min) and Adaboost.R2_SVR (85.34 min) ran
346 slowly among ML methods. Adaboost.R2 based on KRR (Adaboost.R2_KRR) was
347 much more time-saving than Adaboost.R2_SVR. Since the MCMC algorithm required
348 more iteration time to reach convergence, BayesHE was the slowest as expected, and it
349 took 226.12 minutes for each cross-validation.

350 **Discussion**

351 Our results elucidated that ssGBLUP performed better than GBLUP in accuracy in all
352 scenarios investigated, which was consistent with previous studies ^[25, 39-41] . It could be
353 explained by the fact that GBLUP utilized phenotypic information only from genotyped
354 individuals, while ssGBLUP simultaneously used information of both genotyped and
355 non-genotyped individuals to construct a genotype-pedigree relationship matrix (H
356 matrix). Since non-genotyped individuals were related to individuals in the validation
357 population on the pedigree, ssGBLUP took advantage of the phenotypic information of
358 the whole population to obtain better prediction results. However, in our research,
359 ssGBLUP only produced slightly higher accuracies for the two reproduction traits, and
360 the improvements were much lower than those obtained by all ML methods. The lower
361 improvement of ssGBLUP may be due to the following reasons. (I) Weak relationship
362 between the non-genotyped reference population and genotyped candidates in the
363 pedigree. In our study, only 143 of the 789 non-genotyped reference population used
364 by ssGBLUP had pedigree information, and only 46 and 40 individuals' sires and dams
365 were included in the 2566 genotyped individuals, indicating that the relationship
366 between non-genotyped reference animals and genotyped candidates was pretty weak,
367 making tiny contribution to the genomic prediction. Li et al.^[40] showed that the
368 improvement of ssGBLUP over GBLUP on accuracy was almost entirely contributed
369 by non-genotyped close relatives of candidates. It can also be observed from Figure S1
370 that the greater the weight of the A matrix, the lower the accuracy, indicating that the
371 information obtained from pedigree is limited, resulting in ssGBLUP not exerting its

372 advantages greatly. (II) The low heritabilities of TNB and NBA. In this study, the
373 heritabilities for the two traits were both 0.12, which was generally consistent with other
374 reports ^[25, 42, 43], therefore, it cannot get enough accuracy from the pedigree information.
375 This also confirmed by other studies, that a certain improvement can be achieved by
376 adding a smaller reference population for traits with medium or high heritability^[2, 44].

377

378 In this study, we investigated the performance of ML methods in genomic prediction,
379 and demonstrated their superiorities compared to classical methods GBLUP, ssGBLUP
380 and Bayesian methods. Generally, the following characteristics of ML methods make it
381 potentially attractive to genomic prediction. (I) Although ML methods generally require
382 moderate fine-tuning of hyper-parameters, and the default hyper-parameters usually do
383 not perform badly ^[33]. According to our results, the average improvement of ML
384 methods after tuning parameters was 14.3% over using the default hyper-parameters,
385 nonetheless, all ML results without tuning hyper-parameters performed better than
386 GBLUP except for RF in TNB, with an improvement from 0.5% to 8.2% (Table S2).
387 (II) ML methods could handle the number of parameters larger than the sample size, it
388 is very efficient in the case with high-density genetic markers for GS ^[45]. (III) ML
389 methods do not make distribution assumptions about the genetic determinism
390 underlying the trait, enabling to capture the possible non-linear relationships between
391 genotype and phenotype in a flexible way ^[45], and it is different from GBLUP and
392 Bayesian methods, which assumes that all marker effects follow the same normal
393 distribution, or have different classes of shrinkage for different SNP effects. In addition,

394 ML methods can take the correlation and interaction of markers into account as well,
395 while linear models based on pedigree and genomic relationships may not provide a
396 sufficient approximation of the genetic signals generated by complex genetic systems
397 [14]. Consequently, when traits are affected by non-additive effects, especially epistasis,
398 ML methods can achieve more accurate predictions [23]. These make ML methods gain
399 large advantage over GBLUP and BayesHE even they only use genotyped animals.

400 Our results showed that ML methods have improved the prediction accuracy of the
401 reproduction traits in Chinese Yorkshire pig population. SVR, KRR, RF and
402 Adaboost.R2 reflected the superiority of the ML methods, with an average
403 improvement of 20.5%, 21.0%, 14.1% and 20.5% respectively over GBLUP. Liang et
404 al. [46] pointed out that the average improvement of SVR on beef cattle reached a
405 staggering 12.7% . An et al. [13] designed a Cosine kernel-based KRR (K_cRR) and
406 reported that the accuracy of K_cRR was improved by 13.1% compared with GBLUP in
407 three traits of Chinese Simmental beef cattle population. Alves et al.[38] reported SVR
408 has the highest genomic prediction ability in the comparison with GBLUP, BLASSO,
409 Bayesian regularized ANN and RF in the genomic prediction on the reproductive traits
410 of Nellore cattle.

411 Currently, many ML methods are available, and their performance varied in different
412 scenarios. It is difficult to pick the optimal ML method for genomic prediction. In this
413 study, we implemented SVR, KRR, RF and Adaboost.R2 in the genomic prediction. On
414 the whole, SVR and KRR performed best, and our findings were consistent with other
415 studies showing SVR and KRR had been widely used in the genomic prediction [13, 18,

416 ^{23, 47]}. In the present study, for SVR and KRR, we used a non-linear kernel function
417 (RBF kernel) to map the original input data to a high-dimensional feature space and
418 then constructed a linear model in the feature space to estimate GEBVs, and finally
419 constructed a nonlinear model. In all scenarios of this study, the prediction accuracy of
420 SVR and KRR were almost equivalent. One explanation is that the main difference
421 between SVR and KRR is that KRR assumes that most features hardly affect the
422 estimation of GEBVs, so the coefficients of a large number of features are as close to
423 zero as possible, and only certain features have a greater impact on GEBV ^[46]. SVR and
424 KRR were therefore respectively chosen as weak learners for Adaboost.R2.
425 However, Adaboost.R2 did not show the advantages of its integration capabilities
426 compared with single learning algorithms (SVR and KRR). It mainly because the
427 currently SVR and KRR are sufficient to exert prediction abilities, which may limit the
428 benefit of using ensemble learning. Besides, owing to the slow tuning process of
429 Adaboost.R2, we did not precisely tune the hyper-parameters in this research, resulting
430 in slightly lower prediction accuracy than SVR and KRR. One alternative strategy for
431 Adaboost.R2 is integrating more learners. Liang et al. ^[48] developed a stacking
432 ensemble learning framework (SELF) that integrated SVR, KRR, and ENET to predict
433 GEBVs and showed excellent performance. Among all ML methods in this study, RF
434 demonstrated low prediction ability and computational efficiency. The prediction
435 accuracy of RF is mainly affected by the number and maximum depth of decision trees
436 ^[46], but in order to weigh the practical application feasibility of RF, it is impractical to
437 precisely tune the number of trees, resulting in not training the most ideal RF model,

438 thus compromising its prediction accuracy.

439 Although ML significantly outperformed GBLUP and Bayesian methods, one problem
440 should be noted is the hyper-parameter optimization. In this study, the average
441 improvement after tuning parameters was 14.3% over without tuning. Since ML models
442 have multiple hyper-parameters and they are generally sensitive to changes in hyper-
443 parameters, it might be time-consuming to perform strict hyper-parameter adjustments
444 in the process of training models to obtain high accuracies. And the optimal hyper-
445 parameter depends on the character of traits, data sets *etc.*. Usually, the effect of the
446 default hyper-parameters did not perform poorly as discussed above, and failure to find
447 suitable hyper-parameters may greatly reduce the prediction effect of ML methods [46].
448 If hyper-parameter automation can be realized during ML operation, it will greatly
449 reduce the time used for hyper-parameter adjustment and greatly increase the
450 application of ML methods in genomic prediction.

451 **Conclusions**

452 In this study, we compared four ML methods with GBLUP, ssGBLUP and BayesHE to
453 explore their efficiency of genomic prediction on reproduction traits in pigs. We
454 compared the prediction accuracy, unbiasedness, MSE, MAE and computation time of
455 different methods through 20 replicates of 5-fold CV. Our results showed that ML
456 methods possess a significant potential to improve genomic prediction over GBLUP,
457 ssGBLUP and BayesHE. ML methods outperformed in all scenarios, they yielded
458 higher accuracy and smaller MSE and MAE. Among ML methods, SVR and KRR
459 performed the best overall, which yielded higher accuracies, lower bias, and higher

460 computing efficiency. Our findings demonstrated that ML methods are more efficient
461 than traditional genomic selection methods, it could be new options for genomic
462 prediction.

463 **List of abbreviations**

GS	genomic selection
GEBV	genomic breeding values estimation
GBLUP	genomic BLUP
ssGBLUP	single-step GBLUP
ML	machine learning
RF	random forest
SVR	support vector regression
KRR	kernel ridge regression
RKHS	Reproducing Kernel Hilbert space
SVM	support vector machine
TNB	total number of piglets born
NBA	number of piglets born alive
A matrix	pedigree relationship matrix
EBV	estimated breeding values
y_c	corrected phenotypes
DR2	the dosage R-squared measure
COR	the genotype correlation
CR	the genotype concordance rate

MAF	minor allele frequency
BayesHE	Bayesian Horseshoe
5-fold CV	five-fold cross validation
MSE	mean squared error
MAE	mean absolute error
K _c RR	cosine kernel-based KRR
SELF	stacking ensemble learning framework

464

465 **Declarations**

466 **Ethics approval and consent to participate**

467 Animal samples used in this study were approved by the Animal Care and Use
468 Committee of China Agricultural University. There was no use of human participants,
469 data or tissues.

470 **Consent for publication**

471 Not applicable

472 **Availability of data and material**

473 The datasets used or analyzed during the present study are available from the
474 corresponding author on reasonable request.

475 **Competing interests**

476 The authors declare that they have no conflict of interest.

477 **Funding**

478 This work was supported by grants from the National Key Research and Development

479 Project (2019YFE0106800), Modern Agriculture Science and Technology Key Project
480 of Hebei Province (19226376D), China Agriculture Research System of MOF and
481 MARA.

482 **Authors' contributions**

483 XDD designed the experiments. XW performed statistical analysis and wrote the
484 manuscript. SLS provided help on BayesHE. XDD revised the manuscript. All authors
485 read and approved the final manuscript.

486 **Acknowledgement**

487 The authors gratefully acknowledge the constructive comments from reviewers.

488

489 **References**

- 490 1. de Roos AP, Schrooten C, Veerkamp RF, van Arendonk JA. Effects of genomic
491 selection on genetic improvement, inbreeding, and merit of young versus proven bulls.
492 J Dairy Sci. 2011;94(3):1559-67.
- 493 2. Hayes BJ, Bowman PJ, Chamberlain AJ, Goddard ME. Invited review: Genomic
494 selection in dairy cattle: progress and challenges. J Dairy Sci. 2009;92(2):433-43.
- 495 3. Heffner EL, Jannink JL, Sorrells ME. Genomic Selection Accuracy using
496 Multifamily Prediction Models in a Wheat Breeding Program. The Plant Genome.
497 2011;4(1).
- 498 4. Schaeffer LR. Strategy for applying genome-wide selection in dairy cattle. J Anim
499 Breed Genet. 2006;123(4).

- 500 5. García-Ruiz A, Cole JB, VanRaden PM, Wiggans GR, Ruiz-López FJ, Van Tassell
501 CP. Changes in genetic selection differentials and generation intervals in US Holstein
502 dairy cattle as a result of genomic selection. *Proc Natl Acad Sci USA*.
503 2016;113(28):E3995-E4004.
- 504 6. VanRaden PM. Efficient methods to compute genomic predictions. *J Dairy Sci*.
505 2008;91(11):4414-23.
- 506 7. Misztal I, Legarra A, Aguilar I. Computing procedures for genetic evaluation
507 including phenotypic, full pedigree, and genomic information. *J Dairy Sci*.
508 2009;92(9):4648-55.
- 509 8. Christensen OF, Lund MS. Genomic prediction when some animals are not
510 genotyped. *Genet Sel Evol*. 2010;42:2.
- 511 9. Meuwissen THE, Hayes BJ, Goddard ME. Prediction of Total Genetic Value Using
512 Genome-Wide Dense Marker Maps. *Genetics*. 2001;157(4):1819-29.
- 513 10. Habier D, Fernando RL, Kizilkaya K, Garrick DJ. Extension of the bayesian
514 alphabet for genomic selection. *BMC Bioinformatics*. 2011;12:186.
- 515 11. Varona L, Legarra A, Toro MA, Vitezica ZG. Non-additive Effects in Genomic
516 Selection. *Front Genet*. 2018;9:78.
- 517 12. Gianola D, Campos G, Gonzalez-Recio O, Long N, Okut H, Rosa G, et al.
518 *Statistical Learning Methods For Genome-based Analysis Of Quantitative Traits*2018.
- 519 13. An B, Liang M, Chang T, Duan X, Gao HJB. KRR: a nonlinear machine
520 learning with a modified genomic similarity matrix improved the genomic prediction
521 efficiency. *Briefings in Bioinformatics*. 2021.

- 522 14. Gianola D, Okut H, Weigel KA, Rosa GJ. Predicting complex quantitative traits
523 with Bayesian neural networks: a case study with Jersey cows and wheat. *BMC Genet.*
524 2011;12:87.
- 525 15. González-Recio O, Rosa GJM, Gianola D. Machine learning methods and
526 predictive ability metrics for genome-wide prediction of complex traits. *Livestock*
527 *Science.* 2014;166:217-31.
- 528 16. Montesinos-Lopez OA, Martin-Vallejo J, Crossa J, Gianola D, Hernandez-Suarez
529 CM, Montesinos-Lopez A, et al. A Benchmarking Between Deep Learning, Support
530 Vector Machine and Bayesian Threshold Best Linear Unbiased Prediction for
531 Predicting Ordinal Traits in Plant Breeding. *G3 (Bethesda).* 2019;9(2):601-18.
- 532 17. Statnikov A, Wang L, Aliferis CF. A comprehensive comparison of random forests
533 and support vector machines for microarray-based cancer classification. *BMC*
534 *Bioinformatics.* 2008;9:319.
- 535 18. Gonzalez-Camacho JM, Ornella L, Perez-Rodriguez P, Gianola D, Dreisigacker S,
536 Crossa J. Applications of Machine Learning Methods to Genomic Selection in Breeding
537 Wheat for Rust Resistance. *Plant Genome.* 2018;11(2).
- 538 19. Ornella L, Perez P, Tapia E, Gonzalez-Camacho JM, Burgueno J, Zhang X, et al.
539 Genomic-enabled prediction with classification algorithms. *Heredity (Edinb).*
540 2014;112(6):616-26.
- 541 20. Noe F, De Fabritiis G, Clementi C. Machine learning for protein folding and
542 dynamics. *Curr Opin Struct Biol.* 2020;60:77-84.

- 543 21. Kojima K, Tadaka S, Katsuoka F, Tamiya G, Yamamoto M, Kinoshita K. A
544 genotype imputation method for de-identified haplotype reference information by using
545 recurrent neural network. *PLoS Comput Biol.* 2020;16(10):e1008207.
- 546 22. Fa R, Cozzetto D, Wan C, Jones DT. Predicting human protein function with multi-
547 task deep neural networks. *PLoS One.* 2018;13(6):e0198216.
- 548 23. Long N, Gianola D, Rosa GJ, Weigel KA. Application of support vector regression
549 to genome-assisted prediction of quantitative traits. *Theor Appl Genet.*
550 2011;123(7):1065-74.
- 551 24. Madsen P, Jensen J, Labouriau R, Christensen O, Sahana G. DMU - A Package for
552 Analyzing Multivariate Mixed Models in quantitative Genetics and Genomics.
553 Canada August 17-22, 2014.
- 554 25. Guo X, Christensen OF, Ostersen T, Wang Y, Lund MS, Su G. Improving genetic
555 evaluation of litter size and piglet mortality for both genotyped and nongenotyped
556 individuals using a single-step method. *J Anim Sci.* 2015;93(2):503-12.
- 557 26. Browning BL, Browning SR. A unified approach to genotype imputation and
558 haplotype-phase inference for large data sets of trios and unrelated individuals. *Am J*
559 *Hum Genet.* 2009;84(2):210-23.
- 560 27. Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-
561 generation PLINK: rising to the challenge of larger and richer datasets. *Gigascience.*
562 2015;4:7.

- 563 28. Forni S, Aguilar I, Misztal I. Different genomic relationship matrices for single-
564 step analysis using phenotypic, pedigree and genomic information. *Genet Sel Evol.*
565 2011;43:1.
- 566 29. Shi S, Li X, Fang L, Liu A, Su G, Zhang Y, et al. Genomic Prediction Using
567 Bayesian Regression Models With Global-Local Prior. *Front Genet.* 2021;12:628205.
- 568 30. Boser BE, Guyon IM, Vapnik VN. A training algorithm for optimal margin
569 classifiers. *Proceedings of the fifth annual workshop on Computational learning*
570 *theory - COLT '92*1992. p. 144-52.
- 571 31. Müller AC, Guido S. *Introduction to machine learning with Python: A guide for*
572 *data scientists.* Sebastopol: CA: O'Reilly Media, Inc; 2017.
- 573 32. Exterkate P, Groenen PJF, Heij C, van Dijk D. Nonlinear forecasting with many
574 predictors using kernel ridge regression. *International Journal of Forecasting.*
575 2016;32(3):736-53.
- 576 33. Breiman L. Random forests. *Machine Learning.* 2001;45(1):5-32.
- 577 34. Zhou Z. *Machine Learning.* In: Press TU, editor. Beijing,China: Tsinghua
578 University
579 Press2016. p. 247-63.
- 580 35. Drucker H, editor *Improving Regressors using Boosting Techniques.* ICML; 1997.
- 581 36. Shrestha DL, Solomatine DP. Experiments with AdaBoost.RT, an improved
582 boosting scheme for regression. *Neural Comput.* 2006;18(7):1678-710.
- 583 37. Steiger JH. Tests for comparing elements of a correlation matrix. *Psychological*
584 *Bulletin.* 1980;87(2):245-51.

- 585 38. Alves AAC, Espigolan R, Bresolin T, Costa RM, Fernandes Junior GA, Ventura
586 RV, et al. Genome-enabled prediction of reproductive traits in Nellore cattle using
587 parametric models and machine learning methods. *Anim Genet.* 2021;52(1):32-46.
- 588 39. Song H, Ye S, Jiang Y, Zhang Z, Zhang Q, Ding X. Using imputation-based whole-
589 genome sequencing data to improve the accuracy of genomic prediction for combined
590 populations in pigs. *Genet Sel Evol.* 2019;51(1):58.
- 591 40. Li X, Wang S, Huang J, Li L, Zhang Q, Ding X. Improving the accuracy of genomic
592 prediction in Chinese Holstein cattle by using one-step blending. *Genet Sel Evol.*
593 2014;46:66.
- 594 41. Su G, Madsen P, Nielsen US, Mantysaari EA, Aamand GP, Christensen OF, et al.
595 Genomic prediction for Nordic Red Cattle using one-step and selection index blending.
596 *J Dairy Sci.* 2012;95(2):909-17.
- 597 42. Song H, Zhang Q, Ding X. The superiority of multi-trait models with genotype-by-
598 environment interactions in a limited number of environments for genomic prediction
599 in pigs. *J Anim Sci Biotechnol.* 2020;11:88.
- 600 43. Song H, Zhang J, Jiang Y, Gao H, Tang S, Mi S, et al. Genomic prediction for
601 growth and reproduction traits in pig using an admixed reference population. *Journal*
602 *of Animal Science.* 2017;95(8).
- 603 44. Goddard ME, Hayes BJ. Mapping genes for complex traits in domestic animals
604 and their use in breeding programmes. *Nat Rev Genet.* 2009;10(6):381-91.

605 45. Piles M, Bergsma R, Gianola D, Gilbert H, Tusell L. Feature Selection Stability
606 and Accuracy of Prediction Models for Genomic Prediction of Residual Feed Intake in
607 Pigs Using Machine Learning. *Front Genet.* 2021;12:611506.

608 46. Liang M, Miao J, Wang X, Chang T, An B, Duan X, et al. Application of ensemble
609 learning to genomic selection in chinese simmental beef cattle. *J Anim Breed Genet.*
610 2021;138(3):291-9.

611 47. Ghafouri-Kesbi F, Rahimi-Mianji G, Honarvar M, Nejati-Javaremi AJAPS.
612 Predictive ability of Random Forests, Boosting, Support Vector Machines and Genomic
613 Best Linear Unbiased Prediction in different scenarios of genomic evaluation. *Animal*
614 *Production Science.* 2016;57(2):229-36.

615 48. Liang M, Chang T, An B, Duan X, Du L, Wang X, et al. A Stacking Ensemble
616 Learning Framework for Genomic Prediction. *Front Genet.* 2021;12:600040.

617

618 **Table 1** Summary of two reproduction traits of Yorkshire pigs

Trait ^a	Number of records	Birth year	Genotyped animals	Mean	SD	Minimum	Maximum	σ^2_a	σ^2_e	h^2 (SE)
TNB	4274	2016-2020	2566	13	3.38	3	24	1.26	8.95	0.12(0.034)
NBA	4274	2016-2020	2566	12	3.13	3	24	0.98	7.13	0.12(0.032)

619 ^a TNB: total number of piglets born; NBA: number of piglets born alive

620 SE: standard error

621

622 **Table 2** Accuracies and unbiasedness of genomic prediction on TNB and NBA from 7 methods in

Method	TNB		NBA	
	Accuracy	Unbiasedness	Accuracy	Unbiasedness
GBLUP	0.248 ^a ±0.026	0.958±0.132	0.208 ^a ±0.025	0.931±0.142
ssGBLUP	0.251 ^a ±0.026	0.901±0.121	0.221 ^{ab} ±0.026	0.844±0.113
BayesHE	0.243 ^a ±0.025	1.015±0.148	0.207 ^a ±0.026	1.009±0.171
SVR	0.295 ^b ±0.025	1.23±0.119	0.254 ^b ±0.023	1.106±0.11
KRR	0.295 ^b ±0.025	1.266±0.125	0.256 ^b ±0.023	1.151±0.113
RF	0.270 ^{ab} ±0.029	1.229±0.152	0.248 ^{ab} ±0.028	1.188±0.147
Adaboost.R2_SVR	0.293 ^b ±0.025	1.363±0.138	0.254 ^b ±0.024	1.256±0.131
Adaboost.R2_KRR	0.292 ^b ±0.025	1.344±0.136	0.258 ^b ±0.024	1.249±0.129

624 The different superscript of accuracy indicates the significant difference by the Hotelling-Williams

625 test.

626

627 **Table 3** MAE and MSE of 7 methods for TNB and NBA as assessed with 20 replicates of 5-fold

628 CV

Method	TNB		NBA	
	MSE ^a	MAE ^b	MSE ^a	MAE ^b
GBLUP	5.259	1.749	4.168	1.606
ssGBLUP	5.26	1.748	3.95	1.532
BayesHE	5.32	1.763	4.023	1.556
SVR	5.129	1.730	3.880	1.521

KRR	5.134	1.731	3.876	1.521
RF	5.212	1.747	3.901	1.527
Adaboost.R2_SVR	5.158	1.739	3.892	1.528
Adaboost.R2_KRR	5.153	1.737	3.883	1.526

629 ^a MSE: mean squared error

630 ^b MAE: mean absolute error

631

632 **Table 4** Average computing time in one each iteration of the 5-fold Cross validation for different

633 genomic prediction methods

Method	TNB	NBA
GBLUP	2min 06s	2min 02s
ssGBLUP	3min 12s	3min 16s
BayesHE	3h 57min 1s	3h 35min 13s
SVR	5min 27s	5min 07s
KRR	1min 04s	1min 16s
RF	50min 38s	56min 16s
Adaboost.R2_(SVR)	1h 35min 13s	1h 15min 28s
Adaboost.R2_(KRR)	5min 03s	5min 16s

634

635 **Figure captions**

636 **Figure 1 Imputation accuracy**

637 Imputation accuracy of GenoBaits Porcine SNP 50K to PorcineSNP50 BeadChip at
638 different minor allele frequency (MAF) intervals (a) and chromosomes (b).

639 DR², the estimated squared correlation between the estimated allele dose and the true
640 allele dose; Genotype concordance rate (CR), the ratio of the correctly imputed
641 genotypes; Genotype correlation (COR), the correlation coefficient between the
642 imputed variants and the true variants.

643

644

645

Figures

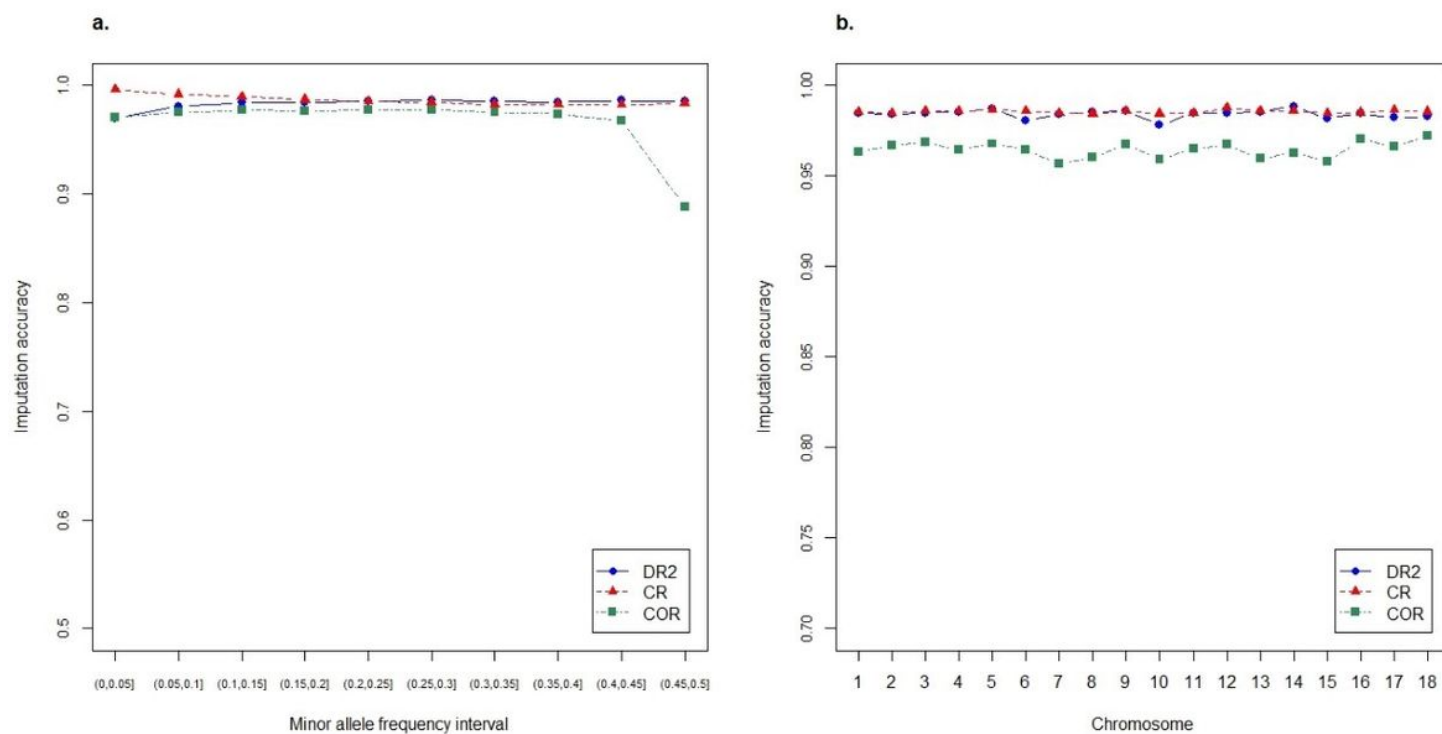


Figure 1

Imputation accuracy of GenoBaits Porcine SNP 50K to PorcineSNP50 BeadChip at different minor allele frequency (MAF) intervals (a) and chromosomes (b). DR2, the estimated squared correlation between the estimated allele dose and the true allele dose; Genotype concordance rate (CR), the ratio of the correctly imputed genotypes; Genotype correlation (COR), the correlation coefficient between the imputed variants and the true variants.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Supplementary.docx](#)