

1 Validation of the linear regression method to evaluate  
2 population accuracy and bias of predictions for  
3 non-linear models

4 Haipeng Yu<sup>1\*</sup>, Rohan L Fernando<sup>2</sup>, and Jack CM Dekkers<sup>2</sup>

5 <sup>1</sup>Department of Animal Sciences, University of Florida, Gainesville, FL,  
6 USA 32611

7 <sup>2</sup>Department of Animal Science, Iowa State University, Ames, IA, USA  
8 50011

9 \* Corresponding author:  
10 Haipeng Yu  
11 Department of Animal Sciences  
12 University of Florida  
13 Gainesville, Florida 32611 USA.  
14 E-mail: haipengyu@ufl.edu  
15

## 16 Abstract

17 Background: The linear regression method (LR) was proposed to estimate population bias and accuracy  
18 of predictions, while addressing the limitations of commonly used cross-validation methods. The validity  
19 and behavior of the LR method have been provided and studied for linear model predictions but not for  
20 non-linear models. The objectives of this study were to 1) provide a mathematical proof for the validity of  
21 the LR method when predictions are based on conditional mean, 2) explore the behavior of the LR method  
22 in estimating bias and accuracy of predictions when the model fitted is different from the true model, and  
23 3) provide guidelines on how to appropriately partition the data into training and validation such that the  
24 LR method can identify presence of bias and accuracy in predictions.

25 Results: We present a mathematical proof for the validity of the LR method to estimate bias and accuracy  
26 of predictions based on the conditional mean, including for non-linear models. Using simulated data, we  
27 show that the LR method can accurately detect bias and estimate accuracy of predictions when an incorrect  
28 model is fitted when the data is partitioned such that the values of relevant predictor variables differ in the  
29 training and validation sets. But the LR method fails when the data are not partitioned in that manner.

30 Conclusions: The LR method was proven to be a valid method to evaluate the population bias and accuracy  
31 of predictions based on the conditional mean, regardless of whether it is a linear or non-linear function of  
32 the data. The ability of the LR method to detect bias and estimate accuracy of predictions when the model  
33 fitted is incorrect depends on how the data are partitioned. To appropriately test the predictive ability of a  
34 model using the LR method, the values of the relevant predictor variables need to be different between the  
35 training and validation sets.

## 36 Background

37 Advances in high-throughput genotyping have enabled the implementation of genomic prediction, which  
38 has facilitated the genetic improvement of animals and plants based on more accurate estimated breeding  
39 values (EBV) at an early stage [e.g., 1–6]. Various genomic prediction models have been proposed and  
40 prediction performance across or within models is usually evaluated by cross-validation (CV) methods [1, 7–  
41 9]. With CV, the data set is partitioned into training and validation sets, with the training set used to fit a  
42 prediction model and estimate the breeding values of individuals in the validation set. Prediction performance  
43 is commonly evaluated with the statistic of predictivity, which is the correlation coefficient between the EBV  
44 and phenotypes adjusted for fixed effects of individuals in the validation set. Scaling predictivity by the  
45 square root of heritability ( $h^2$ ) provides an estimator for prediction accuracy of the EBV [10], defined as the  
46 correlation between true and estimated breeding values. While accuracy estimated with CV has been widely  
47 used to quantify the performance of genomic prediction models, pre-correcting phenotypes in the validation  
48 set using estimates of fixed effects obtained using the whole data set will overestimate the accuracy when  
49 multiple levels of fixed effects are present [11]. Additional limitations include that it can not be applied to  
50 complex models (e.g., random regression models), indirect traits (e.g., unobserved latent traits), and traits  
51 with low heritability ( $h^2$ ) [11].

52 To address these limitations of the CV methodology, Legarra and Reverter [11] proposed a linear regression  
53 (LR) method to estimate the accuracy of genomic prediction. The LR method quantifies the population  
54 accuracy and bias of predictions based on the comparison of EBV of individuals in the validation set estimated  
55 using the training data set with the EBV of those same individuals estimated using the combined training  
56 and validation sets. In the LR method literature, the training set is referred to as the partial data set ( $p$ )  
57 and the combined training and validation data set is referred to as the whole data set ( $w$ ). The LR method  
58 was mathematically proven to provide unbiased estimates of the accuracy and bias of predictions for best  
59 linear unbiased prediction (BLUP) by Legarra and Reverter [11] based on results from Reverter et al. [12].  
60 Macedo et al. [13] investigated the behavior and properties of the LR method by analyzing simulated data  
61 with pedigree-based genetic models. They studied the LR estimators of population bias and accuracy of  
62 predictions by using wrong values of  $h^2$  in the analysis and by fitting wrong models, and claimed that “the

63 LR method works reasonably well for detection of bias when the model used is robust or close to the true  
64 model, and that it works well for estimation of accuracy even when the model is not good". The validity  
65 and performance of the LR method for a non-linear model was explored by Bermann et al. [14]. In their  
66 study, they evaluated the performance of the LR method by fitting a threshold model to both simulation  
67 and real data sets and concluded the LR method can be useful to estimate the directions of bias, dispersion,  
68 and accuracy, though with different magnitudes. The original proof of the LR method [11] was based on the  
69 setting where the whole data set had additional phenotype records relative to the partial data set. Belay et  
70 al. [15] have recently shown that the LR method can also be applied to the setting where the whole data  
71 set has additional genotypes (rather than phenotypes) relative to the partial data set. They used the LR  
72 method to evaluate the bias and accuracy in single-step genomic predictions.

73 While the validity and performance of LR method has been explored using linear and non-linear models  
74 in previous studies [13, 14], a mathematical proof of its validity for non-linear methods of prediction has  
75 not yet been presented. In addition, studies about the performance of the LR method when a model other  
76 than the true model is fitted are still relatively scarce in the literature. The objectives of this study are to  
77 1) present a mathematical proof of the validity of the LR method when predictions are based on conditional  
78 mean, regardless of whether it is a linear or non-linear function of the data 2) investigate the ability of the  
79 LR method to estimate the bias and accuracy of predictions when the fitted model differs from that used to  
80 generate the data, and 3) provide some guidelines on how to partition the data set such that the LR method  
81 can detect bias and estimate accuracy of predictions when the incorrect model is fitted.

## 82 Theory

### 83 Proof that $\text{Cov}(\hat{u}_w, \hat{u}_p) = \text{Var}(\hat{u}_p)$

In the LR method, Legarra and Reverter [11] used  $\text{var}(\mathbf{x})$  to denote the variance of a random element,  $x$ , sampled from a single realization of the random vector  $(\mathbf{x})$ . Here we will denote this variance by  $\text{Var}(x) = \text{var}(\mathbf{x})$ . Let  $u$  denote the breeding value of a validation animal, and  $\hat{u}_p$  and  $\hat{u}_w$  denote the estimated breeding value of  $u$  obtained from partial data and whole data, respectively. Legarra and Reverter [11] proposed the LR method for BLUP by showing  $\text{Cov}(\hat{u}_w, \hat{u}_p) = \text{Var}(\hat{u}_p)$  using the results from Reverter et al. [12] and assumptions of  $\text{Cov}(u, \hat{u}) = \text{Var}(\hat{u})$  and  $\text{E}(\hat{u}_p) = \text{E}(\hat{u}_w) = \text{E}(u)$ . In the following, we prove the validity of the LR method for non-linear models by generalizing the proof of  $\text{Cov}(\hat{u}_w, \hat{u}_p) = \text{Var}(\hat{u}_p)$  for prediction using the conditional mean, which may be non-linear. Let

$$\mathbf{y}_w = \begin{bmatrix} \mathbf{y}_p \\ \mathbf{y}_r \end{bmatrix},$$

where  $\mathbf{y}_w$ ,  $\mathbf{y}_p$ , and  $\mathbf{y}_r$  indicate a vector of phenotype records in the whole, partial, and validation data set, respectively. It is convenient to first show that  $E_{\mathbf{y}_r|\mathbf{y}_p}(\hat{u}_w|\mathbf{y}_p) = \hat{u}_p$ :

$$\begin{aligned} E_{\mathbf{y}_r|\mathbf{y}_p}(\hat{u}_w|\mathbf{y}_p) &= \int_{\mathbf{y}_r} \hat{u}_w \cdot f(\mathbf{y}_r|\mathbf{y}_p) d\mathbf{y}_r \\ &= \int_{\mathbf{y}_r} \int_u u \cdot f(u|\mathbf{y}_p, \mathbf{y}_r) du \cdot f(\mathbf{y}_r|\mathbf{y}_p) d\mathbf{y}_r \\ &= \int_{\mathbf{y}_r} \int_u u \cdot \frac{f(u, \mathbf{y}_p, \mathbf{y}_r)}{f(\mathbf{y}_p, \mathbf{y}_r)} du \cdot f(\mathbf{y}_r|\mathbf{y}_p) d\mathbf{y}_r \\ &= \int_{\mathbf{y}_r} \int_u u \cdot \frac{f(u, \mathbf{y}_p, \mathbf{y}_r)}{f(\mathbf{y}_r|\mathbf{y}_p) \cdot f(\mathbf{y}_p)} du \cdot f(\mathbf{y}_r|\mathbf{y}_p) d\mathbf{y}_r \\ &= \int_u \int_{\mathbf{y}_r} u \cdot \frac{f(u, \mathbf{y}_p, \mathbf{y}_r)}{f(\mathbf{y}_r|\mathbf{y}_p) \cdot f(\mathbf{y}_p)} \cdot f(\mathbf{y}_r|\mathbf{y}_p) d\mathbf{y}_r du \\ &= \int_u u \cdot \frac{f(u, \mathbf{y}_p)}{f(\mathbf{y}_p)} du \\ &= \int_u u \cdot f(u|\mathbf{y}_p) du \\ &= \hat{u}_p \end{aligned}$$

Now, we write the  $Cov(\hat{u}_w, \hat{u}_p)$  as:

$$\begin{aligned} Cov(\hat{u}_w, \hat{u}_p) &= E_{\mathbf{y}_w}[\hat{u}_w(\hat{u}_p - \theta)] \\ &= E_{\mathbf{y}_w}[(\hat{u}_w - \hat{u}_p + \hat{u}_p)(\hat{u}_p - \theta)] \\ &= E_{\mathbf{y}_w}[(\hat{u}_w - \hat{u}_p)(\hat{u}_p - \theta) + \hat{u}_p(\hat{u}_p - \theta)], \end{aligned}$$

where  $\theta$  is the expected value of  $\hat{u}_p$ . But the first term of this expectation can be shown to be null:

$$\begin{aligned} E_{\mathbf{y}_w}[(\hat{u}_w - \hat{u}_p)(\hat{u}_p - \theta)] &= E_{\mathbf{y}_p}\{E_{\mathbf{y}_r|\mathbf{y}_p}[(\hat{u}_w - \hat{u}_p)(\hat{u}_p - \theta)|\mathbf{y}_p]\} \\ &= E_{\mathbf{y}_p}\{(\hat{u}_p - \theta)E_{\mathbf{y}_r|\mathbf{y}_p}[(\hat{u}_w - \hat{u}_p)|\mathbf{y}_p]\} \\ &= E_{\mathbf{y}_p}[(\hat{u}_p - \theta)(\hat{u}_p - \hat{u}_p)] \\ &= 0, \end{aligned}$$

because, as shown previously,  $E_{\mathbf{y}_r|\mathbf{y}_p}(\hat{u}_w|\mathbf{y}_p) = \hat{u}_p$ . Thus, the  $Cov(\hat{u}_w, \hat{u}_p)$  becomes:

$$\begin{aligned} Cov(\hat{u}_w, \hat{u}_p) &= E_{\mathbf{y}_w}[\hat{u}_p \cdot (\hat{u}_p - \theta)] \\ &= E_{\mathbf{y}_p}[\hat{u}_p \cdot (\hat{u}_p - \theta)] \\ &= Var(\hat{u}_p). \end{aligned}$$

84 With the proof of  $Cov(\hat{u}_w, \hat{u}_p) = Var(\hat{u}_p)$ , we showed the LR method holds for non-linear models. This proof  
 85 is similar in principle to that provided by Belay et al. [15], but we recognize that it is not limited to BLUP,  
 86 as invoked in that study, but is applicable to any method of prediction based on the conditional mean [16],  
 87 including for non-linear models.

## 88 Data simulation

A longitudinal data set of body weights in pigs was simulated to evaluate the behavior of LR method for non-linear models, both when the true and a wrong model are used for analysis. Body weights of 1500 individuals from 70 to 500 days of age were simulated using a combination of multi-trait QTL effects (30 bi-allelic QTL) and a Gompertz growth model. Following van Milgen et al. [17], the body weight of individual  $i$  at age  $t$  ( $BW_{it}$ ) was simulated as:

$$BW_{it} = g(\boldsymbol{\theta}_i, t) + \epsilon_{it}, \quad (1)$$

where  $\boldsymbol{\theta}_i = [Age115_i \quad Shape_i \quad BW65_i]$  refers to three underlying latent variables for pig  $i$  of age at 115 kg, a shape parameter, and body weight at 65 days, and  $\epsilon_{it}$  is the residual. We simulated a heterogeneous residuals to mimic the real growth data for pigs using three different residuals across days 70 to 500 ( i.e., 70-167:  $\sigma_{\epsilon_1}^2 = 3.0$ , 168-334:  $\sigma_{\epsilon_2}^2 = 4.0$ , and 335-500:  $\sigma_{\epsilon_3}^2 = 8.0$ ). In equation (1),  $g(\cdot)$  indicates the nonlinear Gompertz function [17]:

$$g(\boldsymbol{\theta}_i, t) = 115 \times \left( \frac{115}{BW65_i} \right)^{\left( -\frac{e^{(-Shape_i(Age115_i - 65))} - e^{(-Shape_i(-65+t))}}{-1 + e^{(-Shape_i(Age115_i - 65))}} \right)}.$$

The three underlying latent variables  $\boldsymbol{\theta}_i$  for individual  $i$  were considered correlated and modeled with a multivariate QTL effects model.

$$\boldsymbol{\theta}_i = \boldsymbol{\mu} + \mathbf{CG}_i + \sum_{j=1}^p m_{ij} \boldsymbol{\alpha}_j + \mathbf{e}_i,$$

89 where  $\boldsymbol{\mu}$  is a vector with the intercepts for each latent variable,  $\mathbf{CG}_i$  is a vector of contemporary group effects,  
 90  $m_{ij}$  is the genotype covariate (0, 1, 2) of individual  $i$  at the  $j$ th QTL,  $\boldsymbol{\alpha}_j$  is a vector of effects for the three  
 91 latent variables for the  $j$ th QTL, and  $\mathbf{e}_i$  is a vector of random environmental effects associated with each  
 92 latent variable. Based on the results of Yu et al. [18], the variance-covariance matrix used for simulation of



93 random environmental effects was equal to 
$$\begin{bmatrix} 3.65 \times 10^{-3} & 1.05 \times 10^{-3} & 5.51 \times 10^{-3} \\ 1.05 \times 10^{-3} & 3.52 \times 10^{-2} & -1.44 \times 10^{-2} \\ -5.51 \times 10^{-3} & -1.44 \times 10^{-2} & 3.05 \times 10^{-2} \end{bmatrix}$$
. The variance-  
 94 covariance used to simulate QTL effects for the three latent variables was arbitrarily but without loss of  
 95 generality derived by dividing the environmental variance-covariance by the number of QTL (i.e., 30).

96 Using the QTL as markers, the simulated data were analyzed with two Bayesian Hierarchical models: 1)  
 97 the Gompertz model that was used for simulation, i.e. the true model, and 2) a quadratic growth model,  
 98 i.e. a wrong model. Variance components that were used to simulate the data were fitted into the true and  
 99 wrong models for analysis. The prediction performances of these two models were evaluated using the LR  
 100 method across 20 replicates. All analyses were performed in Julia [19].

## 101 Data analysis models

We analyzed the simulated data using the Bayesian Hierarchical Gompertz growth model (BHGGM) devel-  
 oped by Yu et al. [18], which integrates a Gompertz growth model, i.e. the true model, with a multi-trait  
 marker effects models. Following equation (1), the three underlying latent variables in the Gompertz growth  
 model were assigned the following prior:

$$\boldsymbol{\theta}_i \sim \mathcal{MVN} \left( \boldsymbol{\mu} + \mathbf{CG}_i + \sum_{j=1}^p m_{ij} \boldsymbol{\alpha}_j, \boldsymbol{\Sigma}_e \right),$$

102 where  $\boldsymbol{\Sigma}_e$  is the environmental variance-covariance matrix, which was assumed to have an inverse Wishart  
 103 prior,  $\mathcal{W}^{-1}(\mathbf{S}_e, \nu_e)$ . The prior for  $\epsilon_{it}$  had a null mean and age specific variances (as described above) to  
 104 allow fitting heterogeneous residuals. Flat priors were assigned to  $\boldsymbol{\mu}$  and  $\mathbf{CG}_i$  and the prior for  $\boldsymbol{\alpha}_j$  followed  
 105  $\mathcal{MVN}(\mathbf{0}, \boldsymbol{\Sigma}_\alpha)$ , where  $\boldsymbol{\Sigma}_\alpha$  has an inverse Wishart distribution,  $\mathcal{W}^{-1}(\mathbf{S}_\alpha, \nu_\alpha)$ .

To fit the Bayesian Hierarchical quadratic growth model (BHQGM), i.e. the wrong model, we introduced  
 a quadratic growth model for the non-linear function  $g(\cdot)$  in equation (1):

$$BW_{it} = b_{i0} + b_{i1}t + b_{i2}t^2 + \epsilon_{it},$$

106 where  $b_{i0}$ ,  $b_{i1}$ , and  $b_{i2}$  refer to three underlying latent variables for individual  $i$  and were assigned the same  
107 multivariate normal prior as  $\theta_i$  in BHGGM. The other parameters in BHQGM also used the same priors as  
108 in BHGGM.

## 109 Design of the partial and validation data sets

110 To investigate the behavior of the LR method, three partitioning scenarios (Figure 1) were implemented: 1)  
111 by animal: phenotype records for days 70 to 500 of the first 500 individuals comprised the partial set and the  
112 phenotypes of the remaining 1000 individuals were assigned into the validation set, 2) by age: phenotypes for  
113 days 70 to 300 of all 1500 individuals comprised the partial set and all 1500 individuals and their phenotypes  
114 from days 301 to 500 were considered as the validation set, and 3) by animal and age: phenotypes for days  
115 70 to 300 of the first 500 individuals comprised the partial set and phenotypes for days 301 to 500 for the  
116 remaining 1000 individuals were assigned to the validation set. The EBV of body weights for individuals in  
117 the validation set across days were then predicted based on the partial set of different scenarios.

## 118 LR method estimators

119 Following Legarra and Reverter [11], estimators of bias, inflation, and accuracy were calculated for EBV of  
120 body weights of the individuals in the validation set using the LR method, as described below.

- 121 1 The LR estimator of population bias is  $\hat{\Delta}_p = \overline{\hat{\mathbf{u}}_p} - \overline{\hat{\mathbf{u}}_w}$ , which is the difference between the mean EBV  
122 of individuals in the validation set estimated from partial and whole data sets, where a value of 0  
123 indicates no bias.
- 124 2 The estimator of inflation is obtained by computing the regression coefficient of  $\hat{u}_w$  on  $\hat{u}_p$ ,  $\hat{b}_{wp} =$   
125  $\frac{\text{Cov}(\hat{u}_w, \hat{u}_p)}{\text{Var}(\hat{u}_p)}$ , which is an estimator of the regression of true breeding value on  $\hat{u}_p$  ( $b_{up}$ ). This estimator  
126 has been referred to as estimator of dispersion by Legarra and Reverter [11] and inflation by Belay  
127 et al. [15], where 1 indicates no inflation. Suppose animals are selected based on EBV to increase the  
128 values of a trait. Then if the true inflation is less than 1, the BV of selected candidates is expected to  
129 be lower than their EBV, which indicates an upward bias of the EBV of the selected animals. On the  
130 other hand, when  $\hat{b}_{wp}$  is larger than 1, the BV of selected candidates will be higher than their EBV,  
131 which indicates an downward bias of the EBV of the selected animals.

132 3 The estimator of population accuracy,  $\hat{\rho}_p = \frac{\text{Cov}(\hat{u}_w, \hat{u}_p)}{\sqrt{\widehat{\text{Var}}(u) \times \text{Var}(\hat{u}_p)}}$ , where  $\widehat{\text{Var}}(u)$  refers to an estimate of the  
133 genetic variance of individuals in the validation set. This estimate was obtained by Gibbs sampling  
134 as:  $\widehat{\text{Var}}(u) = \frac{1}{n_{trn}} \sum_{j=1}^{n_{trn}} u_j^2 - \left( \frac{1}{n_{trn}} \sum_{j=1}^{n_{trn}} u_j \right)^2$ , where  $u_j$  refers to the sampled breeding value of  
135 individual  $j$  and  $n_{trn}$  is the total number of individuals in the training set.

136 In addition to these LR estimators of bias, inflation, and accuracy for body weights predictions, we also  
137 calculated the “true” estimators of these parameters using the simulated values of  $u$  in place of  $\hat{u}_w$  for each  
138 day. Note that these “true” estimators can only be computed in a simulation study, and they are used to  
139 study the performance of the LR estimators, which can be computed in real data analyses.

140 For the true and estimated bias statistics, we calculated their means for each day of age across all animals  
141 in the validation set. These bias statistics were averaged across days within each replicate to test whether  
142 their mean was significantly different from 0 using a t test. Similarly, true and estimated regression coefficient  
143 statistics were averaged across days within each replicate to test whether their mean was significantly different  
144 from 1 using a t test.

## 145 Results

146 To better visualize the prediction performances across the fitted models and partitioning scenarios, we  
147 randomly picked one individual from the validation set and displayed its simulated data against its predictions  
148 in Figure 2. Both simulated body weight phenotypes, true breeding values, and estimated breeding values of  
149 the selected individual were displayed. The predicted data included the breeding values estimated from the  
150 partial and whole data sets for the three partitioning scenarios (Figure 2).

## 151 Population bias

152 Figure 3 shows the true and LR estimates of prediction bias of EBV for body weight at each day when the  
153 data were partitioned by animal. When the true model (BHGGM) was used, both the true and LR estimates  
154 of prediction bias were symmetrically distributed around 0 for each day, and their mean was not significantly  
155 different from 0 ( $P = 0.84$  and  $P = 0.37$ ). In contrast, when the wrong model (BHQGM) was used, the mean  
156 of the true estimates of bias was significantly different from 0 ( $P < 0.001$ ), but the LR estimates of bias  
157 were symmetrically distributed around 0 for each day and their mean was not significantly different from 0  
158 ( $P = 0.4$ ).

159 Figure 4 shows the true and LR estimates of prediction bias of EBV for body weights at each day when  
160 the data were partitioned by age. When the true model was used, both the true and LR estimates of bias  
161 were symmetrically distributed around 0 for each day, and their mean was not significantly different from  
162 0 ( $P = 0.10$  and  $P = 0.09$ ). When the wrong model was used, the true and LR estimates of bias were  
163 significantly different from 0 ( $P < 0.001$  and  $P = 0.002$ ). Results for the partitioning by animal & age were  
164 consistent with those in Figure 4 and are shown in Supplemental Figure S1.

165 Figure 5 shows the true and LR estimates of regression coefficient of EBV for body weights at each day when  
166 the data were partitioned by animal. When the true model was used, both the true and estimated regression  
167 coefficients were symmetrically distributed around 1 for each day, and their mean was not significantly  
168 different from 1 ( $P = 0.75$  and  $P = 0.53$ ). When the wrong model was used, the true and LR estimates of  
169 regression coefficient were significantly different from 1 ( $P < 0.001$ ). Results for the partitioning by age and

170 by animal & age were consistent with those in Figure 5 and are given in Figure 6 and Supplemental Figure  
171 S2, respectively.

## 172 **Population accuracy**

173 In Figure 7, the true and LR estimates of prediction accuracy of EBV for body weights at each day when the  
174 data were partitioned by animal are presented. The LR estimate of prediction accuracy had a similar pattern  
175 as the true estimate of accuracy when using the true model but not when the wrong model was used. When  
176 partitioning the data by age, the LR estimate of accuracy showed a similar pattern as the true estimate of  
177 accuracy curve regardless of which model was fitted. We also evaluated the difference between  $\text{cov}(u, \hat{u}_p)$   
178 and  $\text{cov}(\hat{u}_w, \hat{u}_p)$  when fitting the true and wrong model for the three data partitioning scenarios (Table 1).  
179 There was a non-significant difference ( $P \geq 0.74$ ) between  $\text{cov}(u, \hat{u}_p)$  and  $\text{cov}(\hat{u}_w, \hat{u}_p)$  when the true model  
180 was fitted, but a significant difference ( $P \leq 0.004$ ) was observed for each scenario when the wrong model  
181 was fitted.

## 182 Discussion

183 Based on the initial idea from Reverter et al. [12], Legarra and Reverter [11] proposed the LR method to  
184 quantify the prediction bias and accuracy of EBV at the population level. They proved the validity of LR  
185 method for EBV from a linear model using standard BLUP theory and applied the LR method to a real  
186 cattle data set [11]. While the LR method has also been applied to EBV from a threshold model [14], a  
187 mathematical proof of its validity for a non-linear method of prediction has not been provided. In this study,  
188 we presented a mathematical proof for the validity of the LR method for predictions based on the conditional  
189 mean [16]. In our proof, we assume the partial data contains a subset of the phenotypes in the whole data.  
190 Belay et al. [15] showed the LR method is also applicable to BLUP when the partial data contains a subset  
191 of the genotypes in the whole data. The proof presented in the current paper is similar in principle to that  
192 provided by Belay et al. [15]. Taken together, these two proofs show that the LR method is applicable  
193 to predictions based on conditional mean, regardless of whether the data are partitioned by genotypes or  
194 phenotypes and regardless of whether the model is linear or non-linear. Strictly, however, the LR method is  
195 only valid if the true model is fitted.

196 Using simulated longitudinal data, we confirmed our proof and investigated its behavior when a wrong  
197 model was used for estimation. Furthermore, we explicitly explored how the strategy for partitioning the data  
198 into training and validation sets affect the ability of the LR method to detect bias and estimate accuracy  
199 with the model fitted is not the true model, using three different data partitioning strategies. Below, we  
200 summarize the implications of the fitted model and data partitioning strategies on the performance of the  
201 LR method, thereby providing guidelines for its use to detect bias and estimate accuracy of predictions when  
202 a wrong model is fitted.

203 When the wrong model (BHQM) was fitted and the data were partitioned by animal, the true estimate of  
204 bias was significant, but the LR estimate was not able to identify this bias (Figure 3). Macedo et al. [13] also  
205 observed that for a certain misspecification of the model, the LR method was not able to correctly detect  
206 and estimate a bias. Figure 5 shows that when the wrong model was used, the true estimate of regression  
207 of  $\hat{u}_w$  on  $\hat{u}_p$  had a significant deviation from 1, and in this case the estimate of the regression coefficient  
208 based on the LR method was also significantly different from 1, although differing in magnitude from the

209 true estimate of regression coefficient. This is also consistent with the results observed by Macedo et al. [13].  
210 The pattern of EBV against age were presented in Figure 2 (left column) for a randomly selected individual.  
211 When the wrong model was fitted, the EBV from the partial and whole data sets deviated more from the  
212 true BV than EBV from the true model did. However, even when the wrong model was used, the EBV  
213 from partial and whole data sets were very similar. This explains why the estimate of bias based on the LR  
214 method was not significant when the wrong model was used, although there was a true bias. Figure 7 shows  
215 that, with the wrong model, the accuracy estimated by the LR method was slightly higher than the true  
216 estimate of accuracy, which is consistent with Macedo et al. [13].

217 When the data were partitioned by age, the LR method was able to correctly detect a bias and inflation  
218 when the wrong model was used (Figures 4 and 6). Figure 2 (middle column) shows the EBV of a randomly  
219 selected individual when the data were partitioned by animal. When the wrong model was fitted, the EBV  
220 estimated from partial and whole data sets both deviated from the true BV but the EBV based on the partial  
221 set was quite different from that estimated from the whole set. This illustrates the significant bias that was  
222 detected by the LR method for this scenario. Results for the partitioning by animal & age (right column in  
223 Figure 2) were similar to those when partitioning by age. As in Macedo et al. [13], even with the use of a  
224 wrong model, the accuracy estimated by the LR method was quite close to the true accuracy (Figure 8).

225 The inconsistent bias estimates obtained with the LR method for different data partitioning strategies  
226 suggests that the LR method captures different aspects of the model for different data partitions. When  
227 the data were partitioned by animal, both the partial and whole data sets included phenotypes over the  
228 range from 70 to 500 days. Thus the fit of the growth curve from the partial and whole data sets were  
229 similar, even for the wrong model, although the fit might deviate from that using the true model. Fitting  
230 the wrong growth model is only incorrect in the relationship between age and body weight within individual  
231 but correctly models relationships between relatives. Thus to appropriately test the predictive ability of the  
232 growth model using the LR method, we needed to predict the body weights of animals that are outside the  
233 observed age range for animals in the training set. When the data are partitioned by age, the partial data set  
234 has only body weights measured at ages up to 300 days, whereas the validation data set has body weights  
235 measured at ages up to 500 days. Thus when we predict the body weights of the individuals in the validation  
236 set based on the fit of the growth model from the partial data set, we are testing the predictive ability of the

237 growth model. In this case, the LR method was able to correctly detect a bias when using the wrong model  
238 (bottom right plots in Figure 4 and Figure 6). In real data analyses, repeated k-fold LR can be used to test  
239 the significance of bias. Or if Bayesian method is employed for the LR analysis, the posterior probability of  
240 bias can be computed from a single partitioning of the data.

241 In general, to properly test the predictive ability of a model with the LR method, we need to use the  
242 model to predict the performance of individuals that have values for the relevant predictor variables or  
243 combination of predictor variables that were not present in the training data. In our simulated data, the  
244 predictor variables included the marker genotypes, as well as age. Let's define the predicted performance of  
245 individual  $i$  as  $\hat{y}_i = f(\mathbf{x}_i; \hat{\boldsymbol{\theta}})$ , where  $f(\cdot)$  is the linear or non-linear function used for prediction,  $\mathbf{x}_i$  is a vector  
246 of predictor variables for individual  $i$ , and  $\hat{\boldsymbol{\theta}}$  is estimates of model parameters. Below we will use genomic  
247 prediction by ridge regression BLUP (RR-BLUP) as an example for illustration. To evaluate the predictive  
248 ability of RR-BLUP, the data are partitioned into training and validation sets. The training set is used to  
249 fit the predictive model  $f(\cdot)$  and to estimate the model parameters  $\boldsymbol{\theta}$  (i.e., marker effects). By plugging the  
250 marker effect estimates  $\hat{\boldsymbol{\theta}}$  and observed marker genotypes  $\mathbf{x}$  into  $f(\cdot)$ , the performance of individuals in the  
251 validation set can be predicted. The predictive ability of the model can then be quantified by comparing  
252 the predicted and observed performances of individuals in the validation set. In RR-BLUP, the relevant  
253 predictor variables are the marker genotypes. Thus the same records cannot be used in both the training  
254 and validation sets. In our example, the LR method was used to determine whether it could detect a bias  
255 when a wrong model was used for analysis of longitudinal body weight data. When predicting longitudinal  
256 body weights, the relevant predictor variable is not the genotype but the age of the animal. When the data  
257 were partitioned by animal, the training (partial) and validation sets included phenotypes for animals with  
258 age ranging from days 70 to 500, the same age range as training data was used for the validation data and,  
259 therefore, the LR method failed. However, when the data were partitioned by age, the model was trained  
260 using phenotypes with ages ranging from days 70 to 300 and it was tested by predicting body weights for  
261 animals with age ranging from days 301 to 500. In this case, the LR method was able to detect a bias when  
262 using the wrong model. This was even true when the same genotypes were used in both the training and  
263 validation sets, because to check if the model used for predicting longitudinal body weights is correct, the  
264 relevant predictor variable is age.



## 265 **Conclusions**

266 In the present study, we provide a mathematical proof for the validity of applying the LR method to  
267 predictions based on the conditional mean, regardless of whether it is a linear or non-linear function of data.  
268 Using simulated data, we observed that the LR method was able to detect bias in predictions when an  
269 incorrect non-linear model was fitted. However, when a wrong model is fitted, testing the predictive ability  
270 of the model using the LR method is only valid if the validation set includes values of relevant predictor  
271 variables that are not present in the training set. To our knowledge, this marks the first study that provides  
272 a mathematical proof of the validity of using LR method to a non-linear method of prediction, and we  
273 provide guidelines on how to partition data such that the LR method can detect bias and estimate accuracy  
274 of predictions when the model fitted is incorrect.

## 275 **Acknowledgements**

## 276 **Funding**

277 This work was funded by USDA National Institute of Food and Agriculture award number 2020-67015-31031.

## 278 **Authors' contributions**

279 HY, JCMD, and RLF conceived the research idea. HY and RLF derived a mathematical proof for the validity  
280 of the LR method for predictions based on conditional mean. HY performed the data analyses and drafted  
281 the manuscript. JCMD and RLF edited the manuscript. All authors read and approved the final manuscript.

## 282 **Ethics declarations**

### 283 **Ethics approval and consent to participate**

284 Not applicable.

### 285 **Consent for publication**

286 Not applicable.

### 287 **Competing interests**

288 The authors declare that they have no competing interests.

## 289 Tables

Table 1: Significance (p-values) of tests for the difference between  $\text{Cov}(u, \hat{u}_p)$  and  $\text{Cov}(\hat{u}_w, \hat{u}_p)$  for the three data partitioning scenarios and the two models.

	By animal	By age	By animal & age
True model	0.74	0.85	0.79
Wrong model	0.004	0.002	0.002

## 290 Figure legends

291 Figure 1: Outline of three data partitioning scenarios to create partial and validation sets for the LR method.

292 Figure 2: Example of simulated phenotypes, true breeding values (BV), and estimated breeding values (EBV)  
293 for body weight by age using the true model (TM) and the wrong model (WM) when the partial data set  
294 (Part) was partitioned using different scenarios.

295 Figure 3: True and LR estimates of bias of EBV of body weights at each day when the true or wrong model  
296 was fitted and when partitioning the data by animal. Grey lines are results of 20 simulation replicates, the  
297 red line is the mean of 20 replicates, and the black line indicates bias = 0. P refers to significance of tests  
298 for the difference between true or LR estimate of bias and 0.

299 Figure 4: True and LR estimates of bias of EBV of body weights at each day when the true or wrong model  
300 was fitted and when partitioning the data by age. Grey lines are results of 20 simulation replicates, the red  
301 line is the mean of 20 replicates, and the black line indicates bias = 0. P refers to significance of tests for  
302 the difference between true or LR estimate of bias and 0.

303 Figure 5: True and LR estimates of regression coefficient of EBV of body weights at each day when the true  
304 or wrong model was fitted and when partitioning the data by animal. Grey lines are results of 20 simulation  
305 replicates, the red line is the mean of 20 replicates, and the black line indicates regression coefficient = 1. P  
306 refers to significance of tests for the difference between true or LR estimate of regression coefficient and 1.

307 Figure 6: True and LR estimates of regression coefficient of EBV of body weights at each day when the true  
308 or wrong model was fitted and when partitioning the data by age. Grey lines are results of 20 simulation  
309 replicates, the red line is the mean of 20 replicates, and the black line indicates regression coefficient = 1. P  
310 refers to significance of tests for the difference between true or LR estimate of regression coefficient and 1.

311 Figure 7: True and LR estimates of accuracy when the true or wrong model was fitted and when partitioning  
312 the data by animal. Grey lines are results of 20 simulation replicates, the red line is the mean of 20 replicates,  
313 and the black line indicates accuracy = 1.

314 Figure 8: True and LR estimates of accuracy when the true or wrong model was fitted and when partitioning  
315 the data by age. Grey lines are results of 20 simulation replicates, the red line is the mean of 20 replicates,  
316 and the black line indicates accuracy = 1.

## 317 Supplemental Figure legends

318 Figure S1: True and LR estimates of bias of EBV of body weights at each day when the true or wrong model  
319 was fitted and when partitioning the data by animal & age. Grey lines are results of 20 simulation replicates,  
320 the red line is the mean of 20 replicates, and the black line indicates bias = 0. P refers to significance of tests  
321 for the difference between true or LR estimate of bias and 0.

322 Figure S2: True and LR estimates of regression coefficient of EBV of body weights at each day when the  
323 true or wrong model was fitted and when partitioning the data by animal & age. Grey lines are results of  
324 20 simulation replicates, the red line is the mean of 20 replicates, and the black line indicates regression  
325 coefficient = 1. P refers to significance of tests for the difference between true or LR estimate of regression  
326 coefficient and 1.

327 Figure S3: True and LR estimates of accuracy of EBV of body weights at each day when the true or wrong  
328 model was fitted and when partitioning the data by animal & age. Grey lines are results of 20 simulation  
329 replicates, the red line is the mean of 20 replicates, and the black line indicates accuracy = 1.

330 **Author details**

331 **References**

- 332 1. Meuwissen, T.H., Hayes, B.J., Goddard, M.: Prediction of total genetic value using genome-wide dense marker maps. *Genetics* **157**(4), 1819–1829  
333 (2001)
- 334 2. Dekkers, J., Hospital, F.: The use of molecular genetics in the improvement of agricultural populations. *Nature Reviews Genetics* **3**(1), 22–32  
335 (2002)
- 336 3. Bernardo, R., Yu, J.: Prospects for genomewide selection for quantitative traits in maize. *Crop Science* **47**(3), 1082–1090 (2007)
- 337 4. Habier, D., Fernando, R.L., Kizilkaya, K., Garrick, D.J.: Extension of the bayesian alphabet for genomic selection. *BMC bioinformatics* **12**(1),  
338 1–12 (2011)
- 339 5. Wolc, A., Stricker, C., Arango, J., Settar, P., Fulton, J.E., O'Sullivan, N.P., Preisinger, R., Habier, D., Fernando, R., Garrick, D.J., *et al.*:  
340 Breeding value prediction for production traits in layer chickens using pedigree or genomic relationships in a reduced animal model. *Genetics*  
341 *Selection Evolution* **43**(1), 1–9 (2011)
- 342 6. Morota, G., Koyama, M., M Rosa, G.J., Weigel, K.A., Gianola, D.: Predicting complex traits using a diffusion kernel on genetic markers with an  
343 application to dairy cattle and wheat data. *Genetics Selection Evolution* **45**(1), 1–15 (2013)
- 344 7. Utz, H.F., Melchinger, A.E., Schön, C.C.: Bias and sampling error of the estimated proportion of genotypic variance explained by quantitative trait  
345 loci determined from experimental data in maize using cross validation and validation with independent samples. *Genetics* **154**(4), 1839–1849  
346 (2000)
- 347 8. Saatchi, M., McClure, M.C., McKay, S.D., Rolf, M.M., Kim, J., Decker, J.E., Taxis, T.M., Chapple, R.H., Ramey, H.R., Northcutt, S.L., *et al.*:  
348 Accuracies of genomic breeding values in american angus beef cattle using k-means clustering for cross-validation. *Genetics Selection Evolution*  
349 **43**(1), 1–16 (2011)
- 350 9. Morota, G., Gianola, D.: Kernel-based whole-genome prediction of complex traits: a review. *Frontiers in genetics* **5**, 363 (2014)
- 351 10. Serão, N.V., Kemp, R.A., Mote, B.E., Willson, P., Harding, J., Bishop, S.C., Plastow, G.S., Dekkers, J.: Genetic and genomic basis of antibody  
352 response to porcine reproductive and respiratory syndrome (prrs) in gilts and sows. *Genetics Selection Evolution* **48**(1), 1–15 (2016)
- 353 11. Legarra, A., Reverter, A.: Semi-parametric estimates of population accuracy and bias of predictions of breeding values and future phenotypes  
354 using the lr method. *Genetics Selection Evolution* **50**(1), 1–18 (2018)
- 355 12. Reverter, A., Golden, B., Bourdon, R., Brinks, J.: Detection of bias in genetic predictions. *Journal of animal science* **72**(1), 34–37 (1994)
- 356 13. Macedo, F., Reverter, A., Legarra, A.: Behavior of the linear regression method to estimate bias and accuracies with correct and incorrect genetic  
357 evaluation models. *Journal of Dairy Science* **103**(1), 529–544 (2020)
- 358 14. Bermann, M., Legarra, A., Hollifield, M.K., Masuda, Y., Lourenco, D., Misztal, I.: Validation of single-step gblup genomic predictions from  
359 threshold models using the linear regression method: An application in chicken mortality. *Journal of Animal Breeding and Genetics* **138**(1), 4–13  
360 (2021)
- 361 15. Belay, T.K., Eikje, L.S., Gjuvsland, A.B., Nordbø, Ø., Tribout, T., Meuwissen, T.: Correcting for base-population differences and unknown parent  
362 groups in single-step genomic predictions of norwegian red cattle. *Journal of Animal Science* (2022)
- 363 16. Fernando, R., Gianola, D.: Optimal properties of the conditional mean as a selection criterion. *Theoretical and applied genetics* **72**(6), 822–825  
364 (1986)
- 365 17. van Milgen, J., Valancogne, A., Dubois, S., Dourmad, J.-Y., Sève, B., Noblet, J.: Inraporc: a model and decision support tool for the nutrition of  
366 growing pigs. *Animal Feed Science and Technology* **143**(1-4), 387–405 (2008)
- 367 18. Yu, H., van Milgen, J., Knol, E.F., Fernando, R.L., Dekkers, J.C.: A bayesian hierarchical model to integrate a mechanistic growth model in  
368 genomic prediction. In: *Proceedings of the World Congress on Genetics Applied to Livestock Production: 3-8 July 2022; Rotterdam, The*  
369 *Netherlands.* (2022)
- 370 19. Bezanson, J., Edelman, A., Karpinski, S., Shah, V.B.: Julia: A fresh approach to numerical computing. *SIAM review* **59**(1), 65–98 (2017)

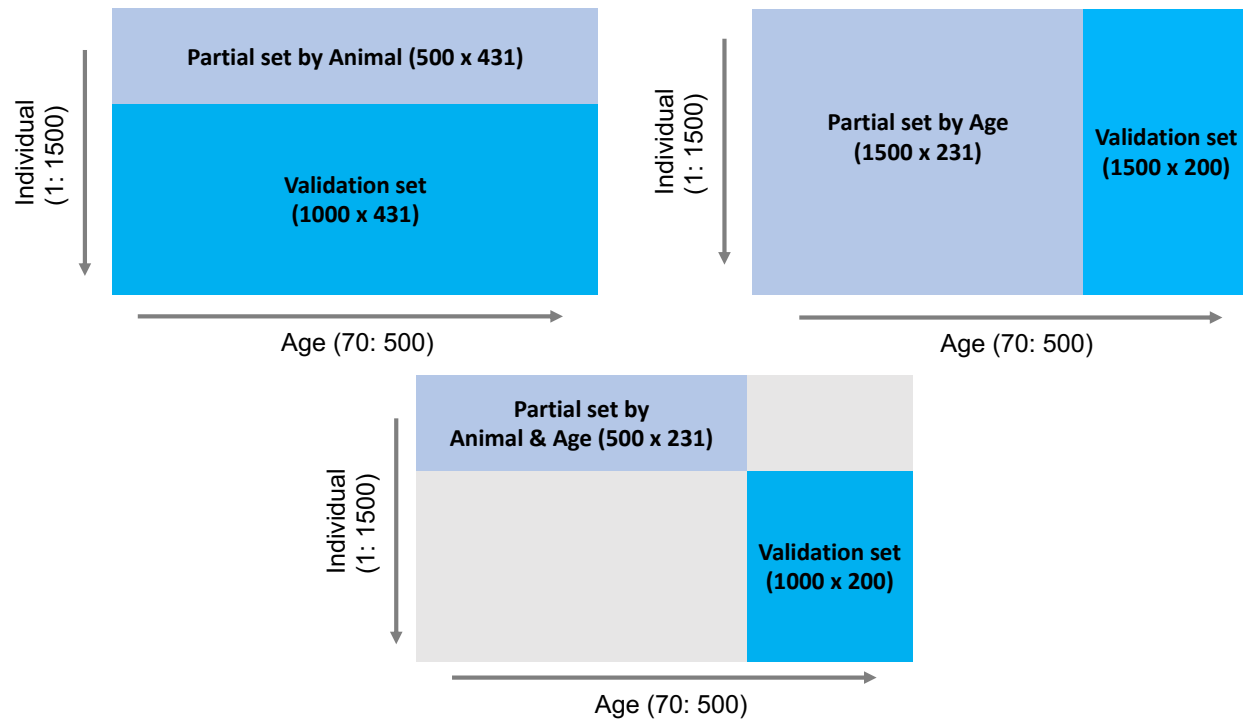


Figure 1: Outline of three data partitioning scenarios to create partial and validation sets for the LR method.



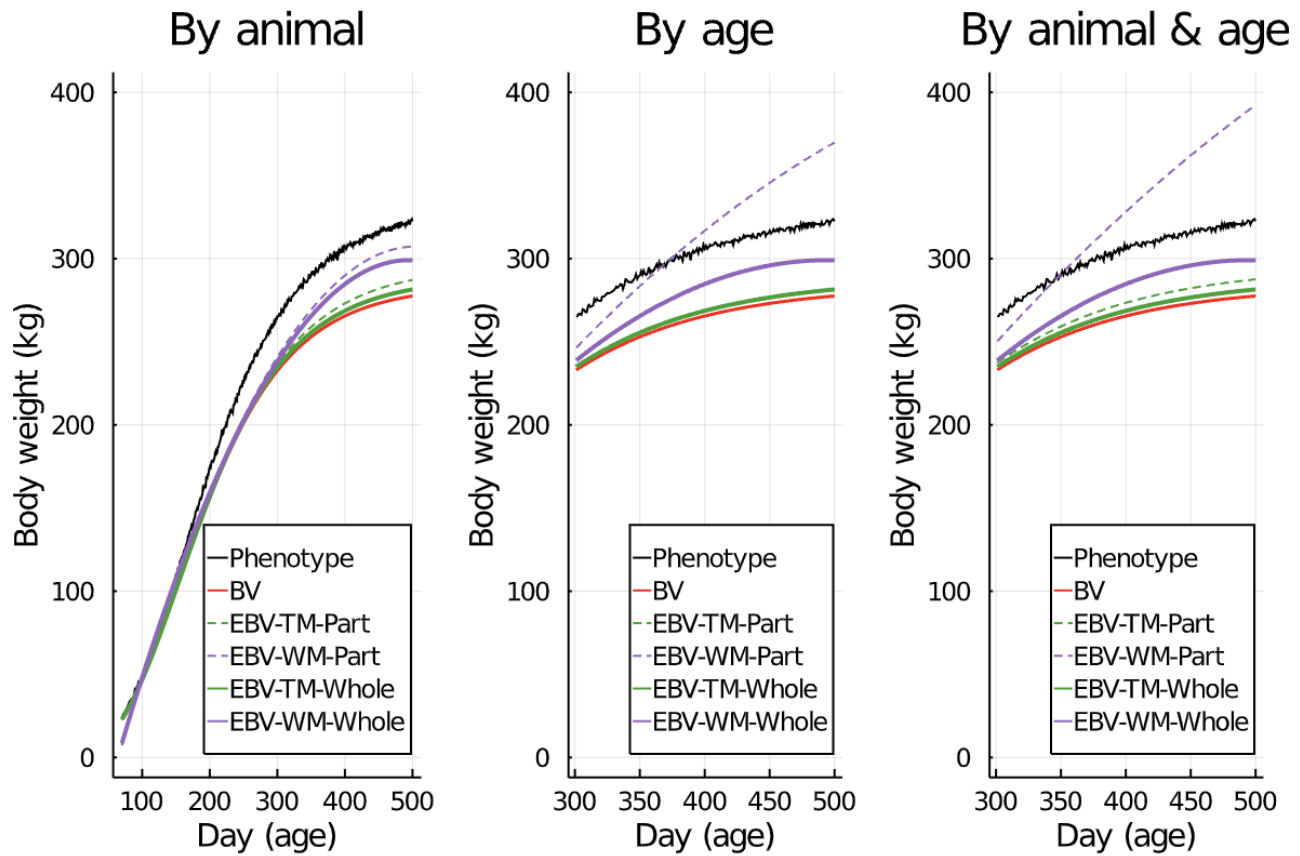


Figure 2: Example of simulated phenotypes, true breeding values (BV), and estimated breeding values (EBV) for body weight by age using the true model (TM) and the wrong model (WM) when the partial data set (Part) was partitioned using different scenarios.

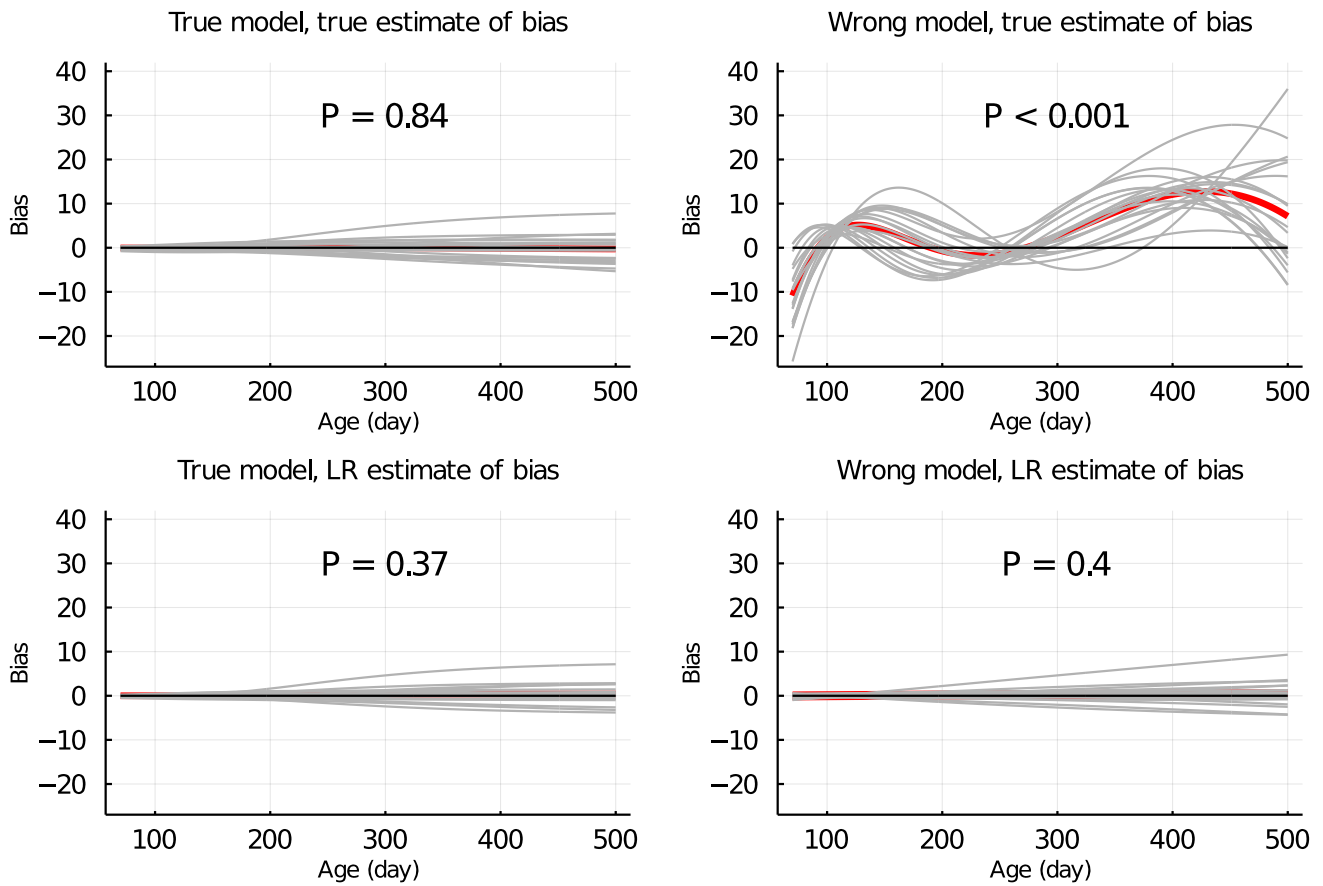


Figure 3: True and LR estimates of bias of EBV of body weights at each day when the true or wrong model was fitted and when partitioning the data by animal. Grey lines are results of 20 simulation replicates, the red line is the mean of 20 replicates, and the black line indicates bias = 0. P refers to significance of tests for the difference between true or LR estimate of bias and 0.

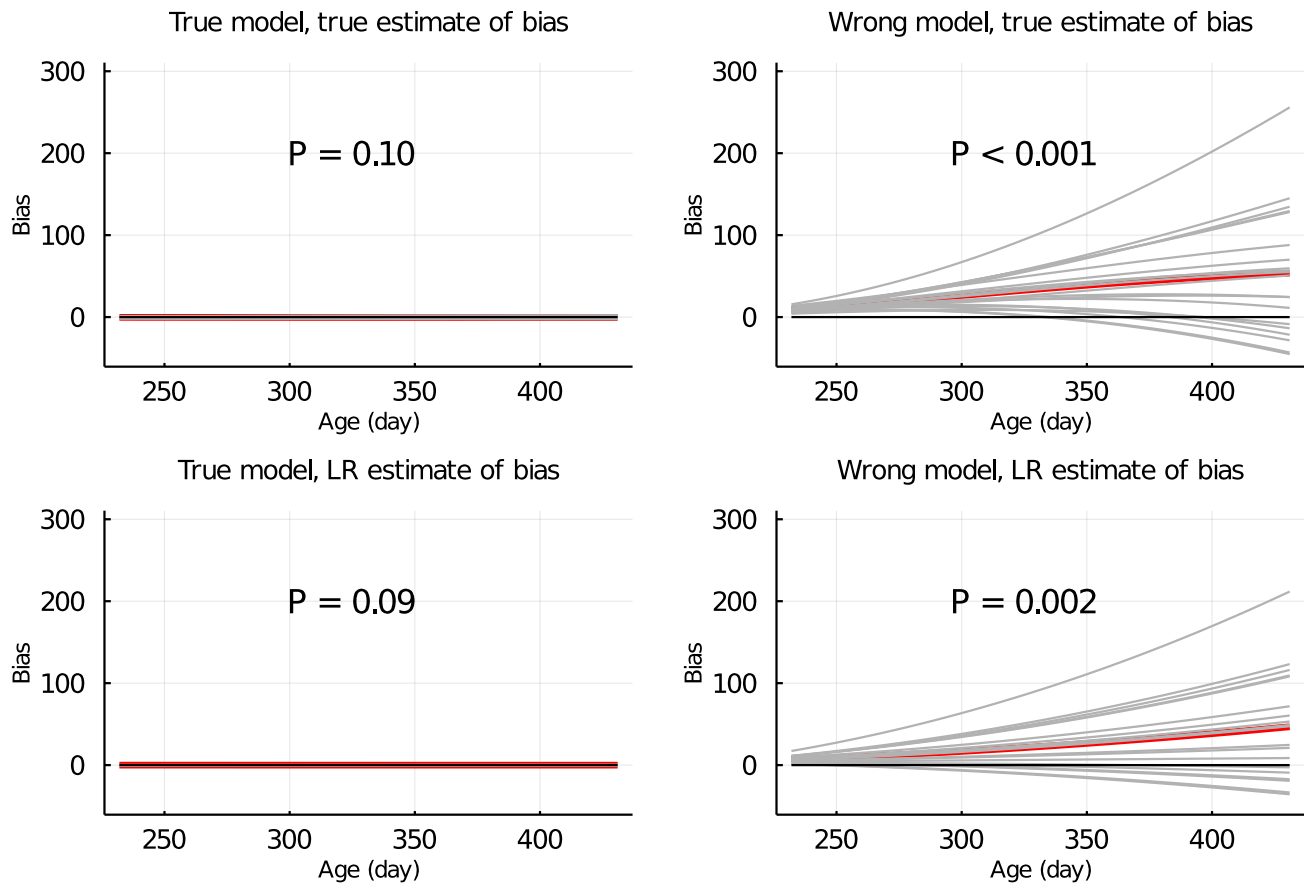


Figure 4: True and LR estimates of bias of EBV of body weights at each day when the true or wrong model was fitted and when partitioning the data by age. Grey lines are results of 20 simulation replicates, the red line is the mean of 20 replicates, and the black line indicates bias = 0. P refers to significance of tests for the difference between true or LR estimate of bias and 0.

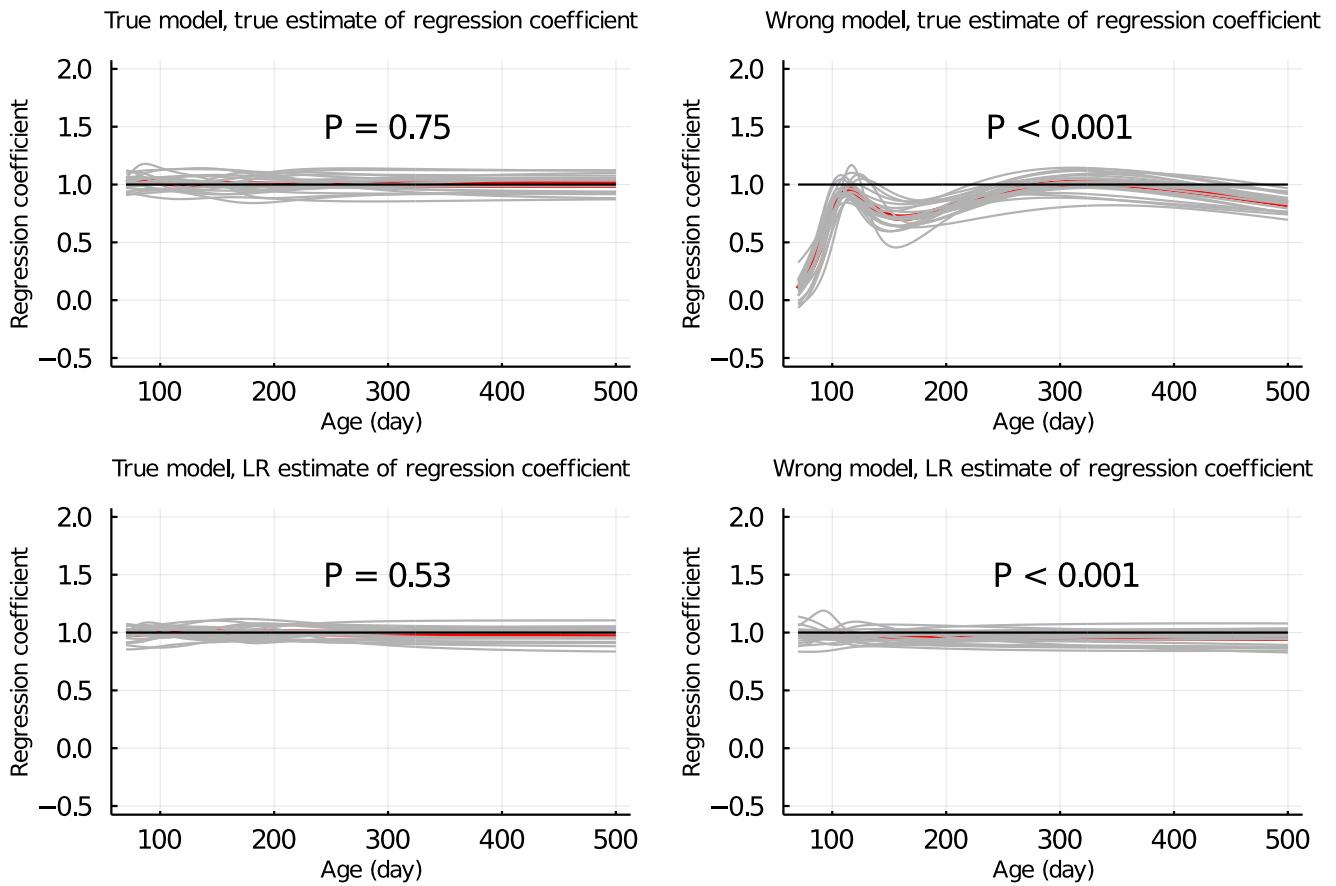


Figure 5: True and LR estimates of regression coefficient of EBV of body weights at each day when the true or wrong model was fitted and when partitioning the data by animal. Grey lines are results of 20 simulation replicates, the red line is the mean of 20 replicates, and the black line indicates regression coefficient = 1. P refers to significance of tests for the difference between true or LR estimate of regression coefficient and 1.

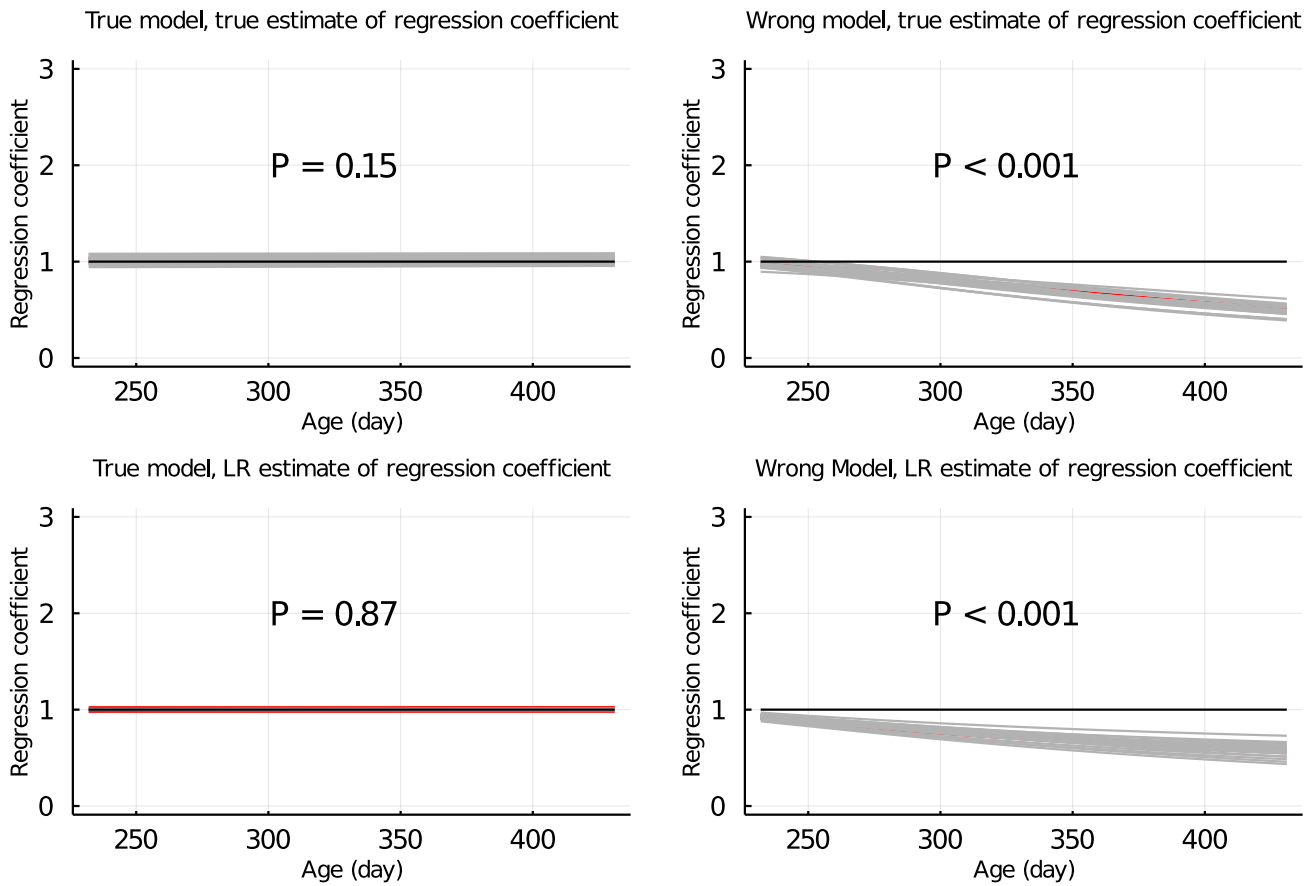


Figure 6: True and LR estimates of regression coefficient of EBV of body weights at each day when the true or wrong model was fitted and when partitioning the data by age. Grey lines are results of 20 simulation replicates, the red line is the mean of 20 replicates, and the black line indicates regression coefficient = 1. P refers to significance of tests for the difference between true or LR estimate of regression coefficient and 1.

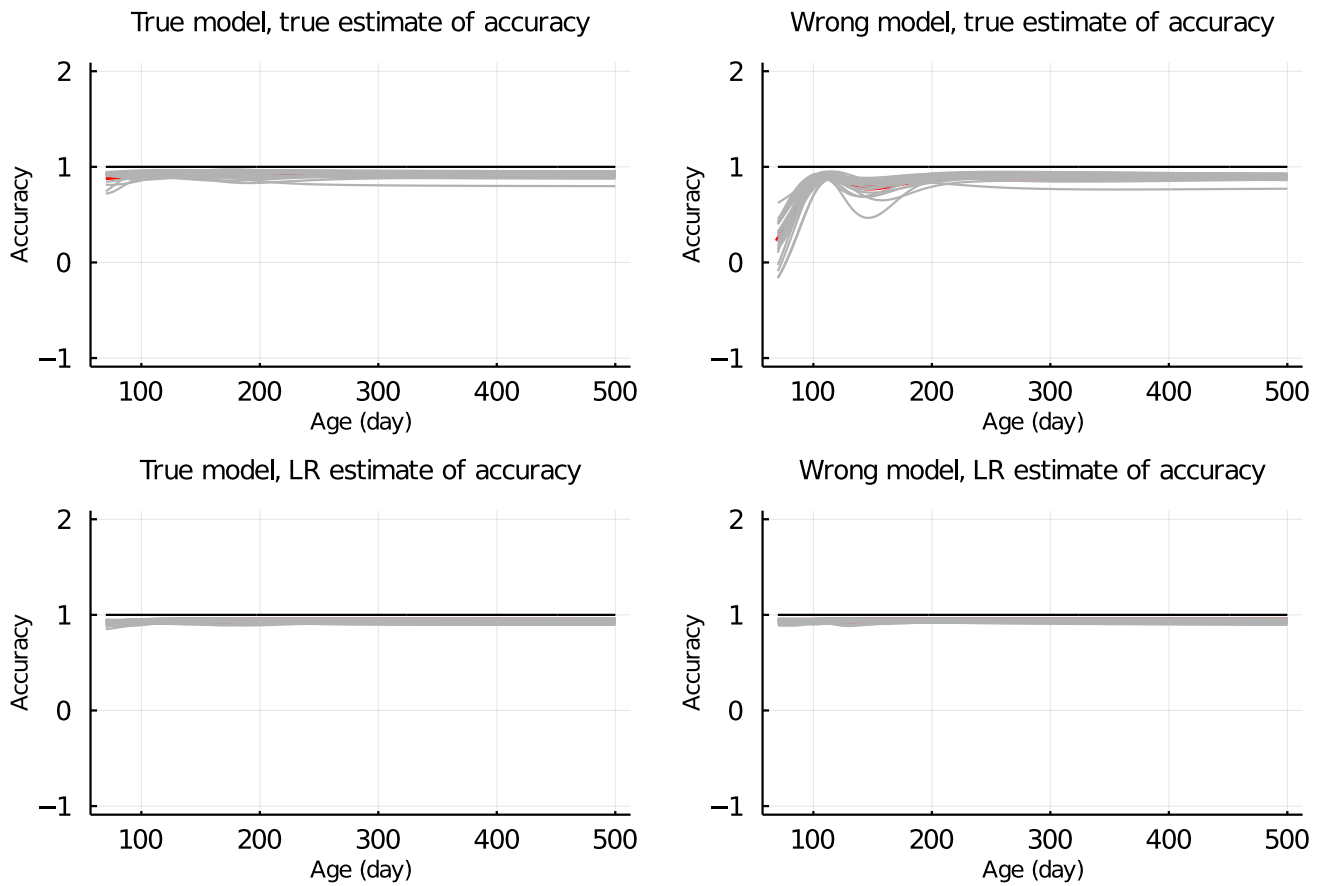


Figure 7: True and LR estimates of accuracy when the true or wrong model was fitted and when partitioning the data by animal. Grey lines are results of 20 simulation replicates, the red line is the mean of 20 replicates, and the black line indicates accuracy = 1.

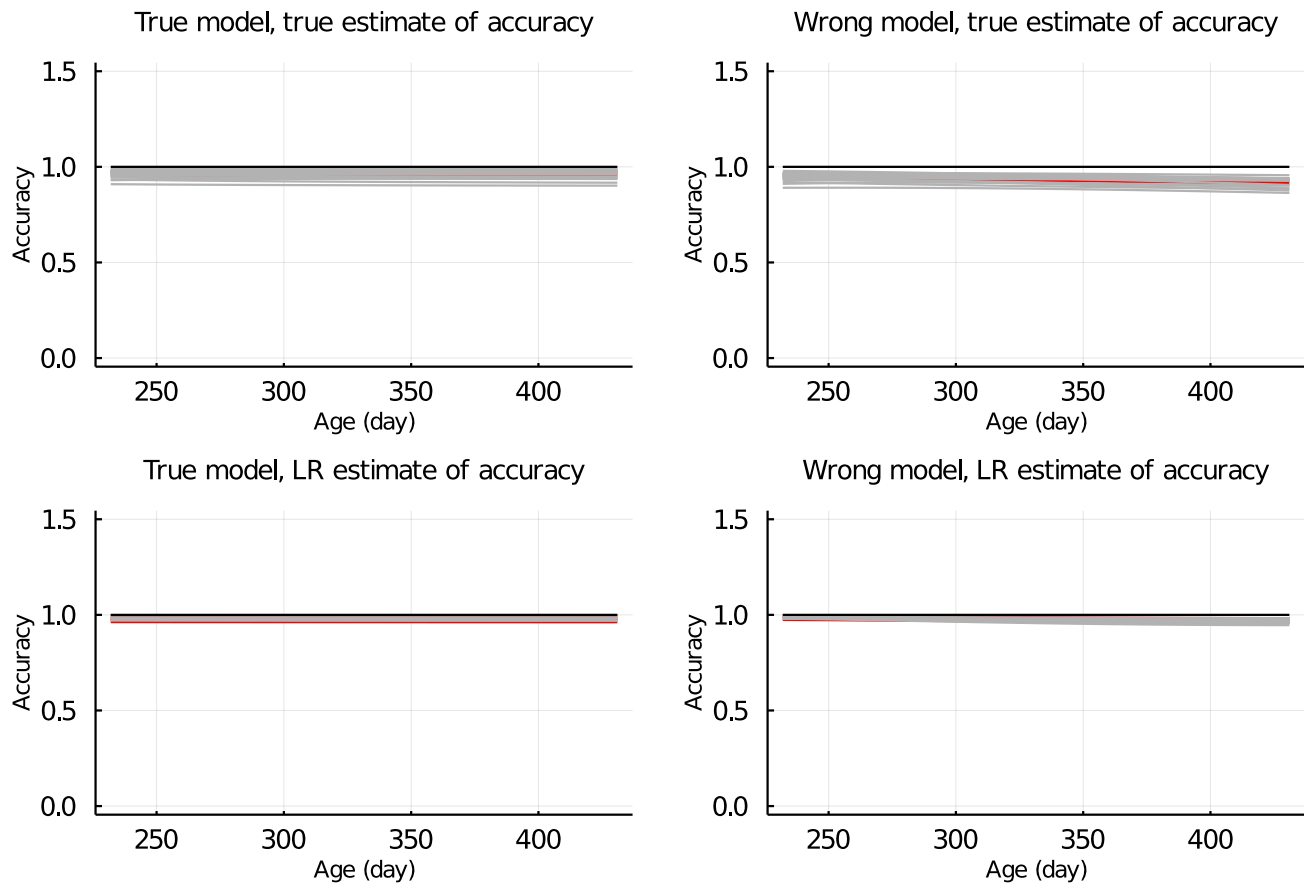


Figure 8: True and LR estimates of accuracy when the true or wrong model was fitted and when partitioning the data by age. Grey lines are results of 20 simulation replicates, the red line is the mean of 20 replicates, and the black line indicates accuracy = 1.

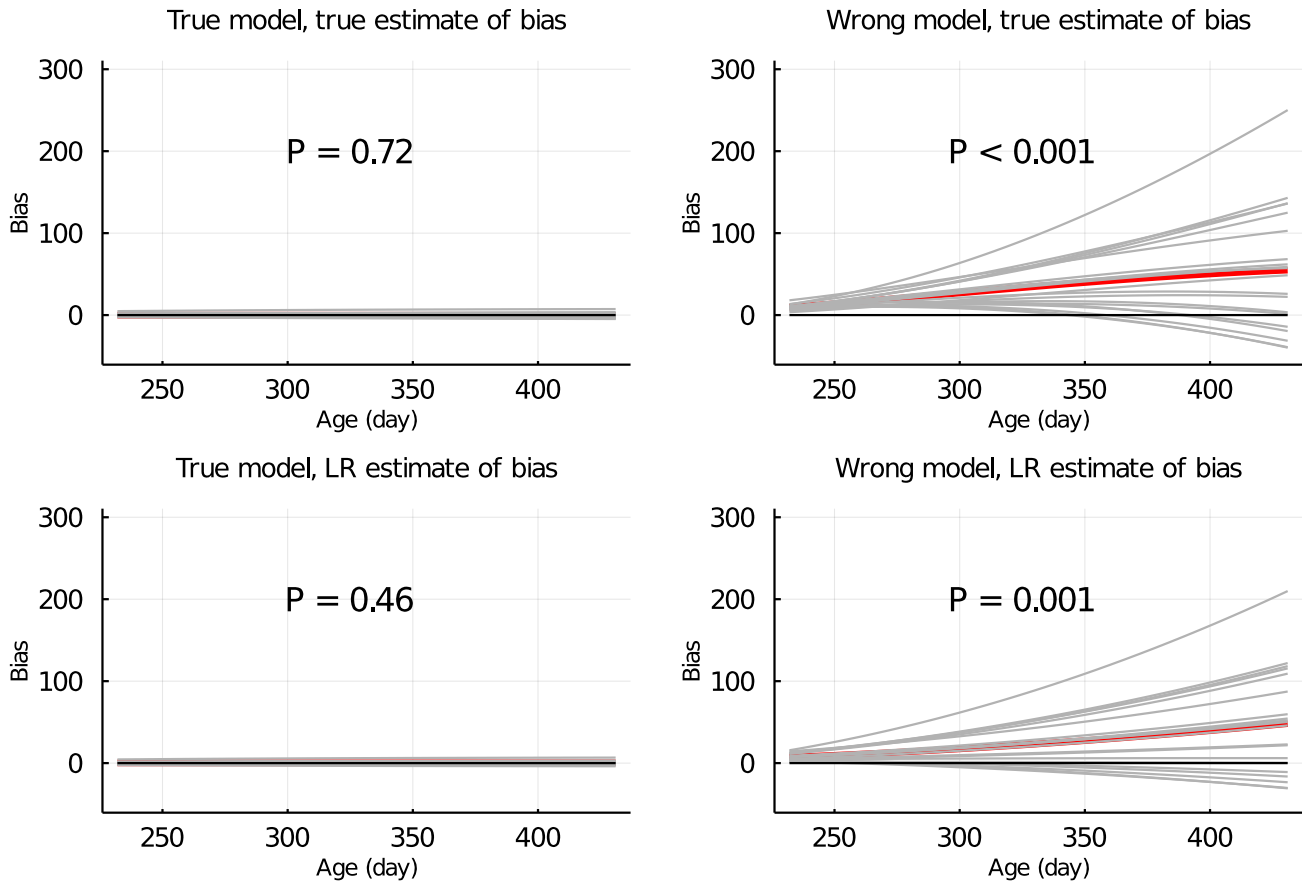


Figure S1: True and LR estimates of bias of EBV of body weights at each day when the true or wrong model was fitted and when partitioning the data by animal & age. Grey lines are results of 20 simulation replicates, the red line is the mean of 20 replicates, and the black line indicates bias = 0. P refers to significance of tests for the difference between true or LR estimate of bias and 0.



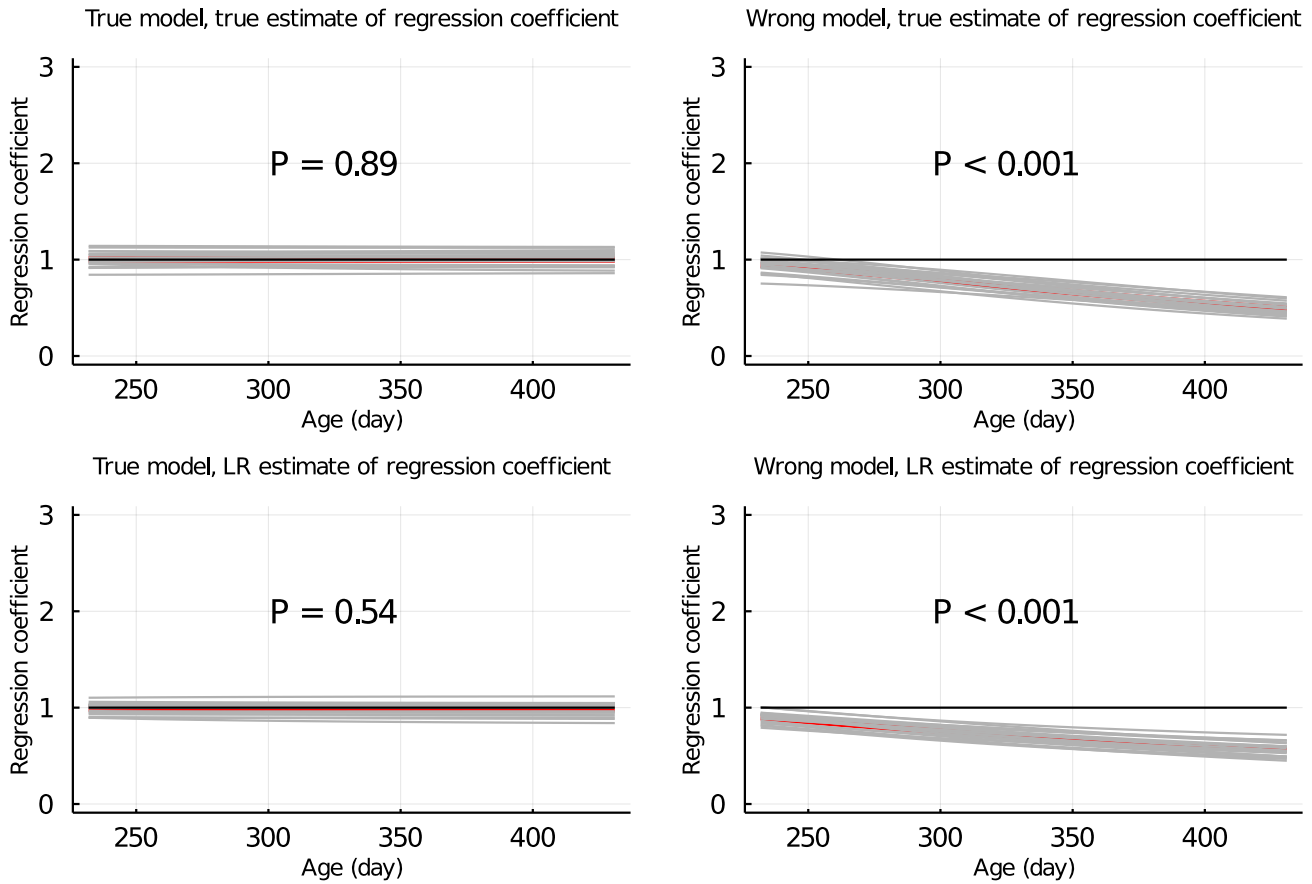


Figure S2: True and LR estimates of regression coefficient of EBV of body weights at each day when the true or wrong model was fitted and when partitioning the data by animal & age. Grey lines are results of 20 simulation replicates, the red line is the mean of 20 replicates, and the black line indicates regression coefficient = 1. P refers to significance of tests for the difference between true or LR estimate of regression coefficient and 1.

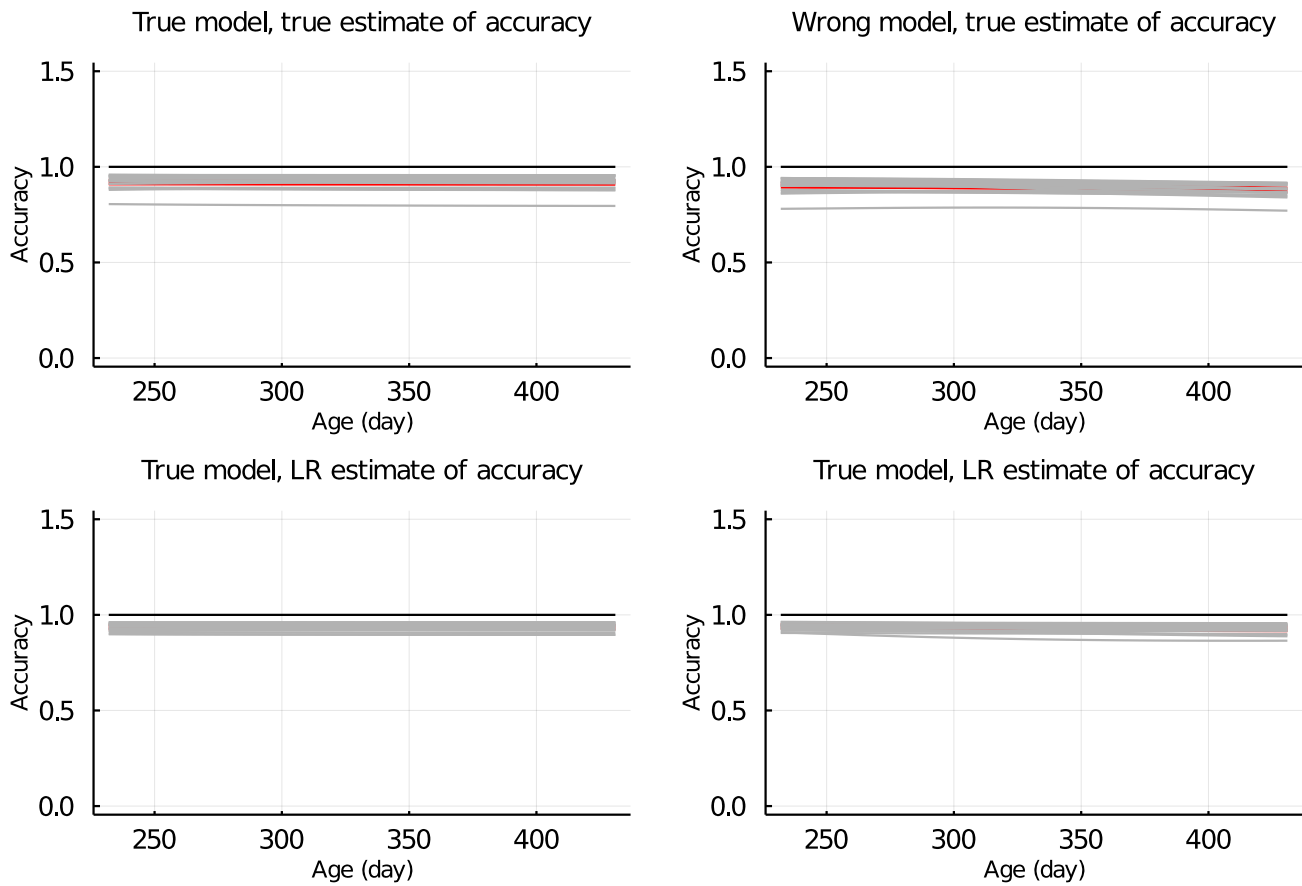


Figure S3: True and LR estimates of accuracy of EBV of body weights at each day when the true or wrong model was fitted and when partitioning the data by animal & age. Grey lines are results of 20 simulation replicates, the red line is the mean of 20 replicates, and the black line indicates accuracy = 1.