

University of Nebraska - Lincoln

DigitalCommons@University of Nebraska - Lincoln

Agronomy & Horticulture -- Faculty Publications

Agronomy and Horticulture Department

2019

Predicting Longitudinal Traits Derived from High-Throughput Phenomics in Contrasting Environments Using Genomic Legendre Polynomials and B-Splines

Mehdi Momen

Virginia Polytechnic Institute and State University

Malachy T. Campbell

Virginia Polytechnic Institute and State University, campbell.malachy@gmail.com

Harkamal Walia

University of Nebraska-Lincoln, hwalia2@unl.edu

Gota Morota

University of Nebraska- Lincoln, morota@vt.edu

Follow this and additional works at: <https://digitalcommons.unl.edu/agronomyfacpub>



Part of the [Agricultural Science Commons](#), [Agriculture Commons](#), [Agronomy and Crop Sciences Commons](#), [Botany Commons](#), [Horticulture Commons](#), [Other Plant Sciences Commons](#), and the [Plant Biology Commons](#)

Momen, Mehdi; Campbell, Malachy T.; Walia, Harkamal; and Morota, Gota, "Predicting Longitudinal Traits Derived from High-Throughput Phenomics in Contrasting Environments Using Genomic Legendre Polynomials and B-Splines" (2019). *Agronomy & Horticulture -- Faculty Publications*. 1289. <https://digitalcommons.unl.edu/agronomyfacpub/1289>

This Article is brought to you for free and open access by the Agronomy and Horticulture Department at DigitalCommons@University of Nebraska - Lincoln. It has been accepted for inclusion in Agronomy & Horticulture -- Faculty Publications by an authorized administrator of DigitalCommons@University of Nebraska - Lincoln.

Predicting Longitudinal Traits Derived from High-Throughput Phenomics in Contrasting Environments Using Genomic Legendre Polynomials and B-Splines

Mehdi Momen,* Malachy T. Campbell,* Harkamal Walia,[†] and Gota Morota*¹

*Department of Animal and Poultry Sciences, Virginia Polytechnic Institute and State University, Blacksburg, VA, 24061 and [†]Department of Agronomy and Horticulture, University of Nebraska-Lincoln, Lincoln, NE, 68583

ORCID IDs: 0000-0002-2562-2741 (M.M.); 0000-0002-8257-3595 (M.T.C.); 0000-0002-9712-5824 (H.W.); 0000-0002-3567-6911 (G.M.)

ABSTRACT Recent advancements in phenomics coupled with increased output from sequencing technologies can create the platform needed to rapidly increase abiotic stress tolerance of crops, which increasingly face productivity challenges due to climate change. In particular, high-throughput phenotyping (HTP) enables researchers to generate large-scale data with temporal resolution. Recently, a random regression model (RRM) was used to model a longitudinal rice projected shoot area (PSA) dataset in an optimal growth environment. However, the utility of RRM is still unknown for phenotypic trajectories obtained from stress environments. Here, we sought to apply RRM to forecast the rice PSA in control and water-limited conditions under various longitudinal cross-validation scenarios. To this end, genomic Legendre polynomials and B-spline basis functions were used to capture PSA trajectories. Prediction accuracy declined slightly for the water-limited plants compared to control plants. Overall, RRM delivered reasonable prediction performance and yielded better prediction than the baseline multi-trait model. The difference between the results obtained using Legendre polynomials and that using B-splines was small; however, the former yielded a higher prediction accuracy. Prediction accuracy for forecasting the last five time points was highest when the entire trajectory from earlier growth stages was used to train the basis functions. Our results suggested that it was possible to decrease phenotyping frequency by only phenotyping every other day in order to reduce costs while minimizing the loss of prediction accuracy. This is the first study showing that RRM could be used to model changes in growth over time under abiotic stress conditions.

KEYWORDS

genomic
prediction
phenomics
longitudinal
modeling
random
regression
time series
GenPred
Shared Data
Resources

BACKGROUND

Plant biology has become a large-scale, data-rich field with the development of high-throughput technologies for genomics and phenomics. This has increased the feasibility of data driven approaches to be applied

to address the challenge of developing climate-resilient crops (Tester and Langridge 2010). Crop responses to environmental changes are highly dynamic and have a strong temporal component. Such responses are also known as function-valued traits, for which means and covariances along the trajectory change continuously. Single time point measurements of phenotypes, however, only provide a snapshot, posing a series of challenges for research efforts aimed at understanding the ability of the plant to mount a tolerant response to an environmental constraint. Advancements in high-throughput phenotyping (HTP) technologies have enabled the automated collection of measurements at high temporal resolution to produce high density image data that can capture a plethora of morphological and physiological measurements (Furbank and Tester 2011). In particular, image-based phenotyping has been deemed a game changer because conventional phenotyping is laborious and often involves destructive methods, precluding

Copyright © 2019 Momen *et al.*

doi: <https://doi.org/10.1534/g3.119.400346>

Manuscript received May 13, 2019; accepted for publication August 15, 2019; published Early Online August 19, 2019.

This is an open-access article distributed under the terms of the Creative Commons Attribution 4.0 International License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

Supplemental material available at FigShare: <https://doi.org/10.25387/g3.9383543>.

¹Corresponding author: Department of Animal and Poultry Sciences, Virginia Polytechnic Institute and State University, 175 West Campus Drive, Blacksburg, Virginia 24061. E-mail: morota@vt.edu

repeated sampling over time from the same individual (Ge *et al.* 2016). More importantly, these HTP systems offer greater potential to uncover the time-specific molecular events driven by important genes that have yet to be discovered in genome-wide association studies (GWAS) or to perform forecasting of future phenotypes in longitudinal genomic prediction. Thus, integrating these HTP data into quantitative genetics has the potential to increase the rate of genetic gain in crops. However, to take full advantage of such opportunities, novel statistical methods that can fully leverage time series HTP data need to be developed.

Recently, Campbell *et al.* (2018) used a random regression model (RRM) to perform genomic prediction for longitudinal HTP traits in controlled environments, such as greenhouses, using Legendre polynomials as the choice of a basis function to model dependencies across time. They also demonstrated that RRM could be used to achieve reasonable prediction accuracy in a cross-validation (CV) framework to forecast future phenotypes based on information from earlier growth stages. RRM also enables the calculation of (co)variances and genetic values at any time between the beginning and end of the trajectory, even including time points that lack phenotypic information. This study showed that RRM could effectively describe the temporal dynamics of genetic signals by accounting for mean and covariance structures that are subjected to change over time (Kirkpatrick *et al.* 1990). However, the utility of RRM for plants under an abiotic stress environment is not explored. This is a critical unknown as the crop productivity is greatly limited by environmental challenges such as drought and heat stress. In addition to the Legendre polynomials, spline functions can be used to describe the relationships between image-based phenomics and genomics data in longitudinal modeling (White *et al.* 1999). In particular, B-spline functions have been used to study a variety of traits, such as growth records, in animal breeding in terms of model goodness of fit using pedigree data (*e.g.*, Meyer 2005; Baldi *et al.* 2010); however, its application to HTP data in plants and its predictive ability from a CV perspective remains untested.

Here we present our results from the performance of RRM applied to HTP temporal shoot biomass data in response to control and water-limited conditions using Legendre polynomials and spline functions. We selected drought stress because water limitation significantly impacts shoot growth and is the major limitation for agricultural productivity and global food security.

MATERIALS AND METHODS

Plant materials and greenhouse conditions

Three hundred fifty-seven accessions ($n = 357$) of the rice (*O. Sativa*) diversity panel 1 (RDP1) were selected for this study (Zhao *et al.* 2011). Seeds were surface sterilized with Thiram fungicide and germinated on moist paper towels in plastic boxes for three days. For each accession, three uniformly germinated seedlings were selected and transplanted to pots (150mm diameter \times 200 mm height) filled with 2.5 kg of UC Mix. Square containers were placed below each pot to allow water to collect. The plants were grown in saturated soil on greenhouse benches prior to phenotyping.

All lines were genotyped with 44,000 single nucleotide polymorphisms (SNPs) (Zhao *et al.* 2011). We used PLINK v1.9 software (Purcell *et al.* 2007) to remove SNPs with a call rate ≤ 0.95 and a minor allele frequency ≤ 0.05 . Missing genotypes were imputed using Beagle software version 3.3.2 (Browning and Browning 2007). Finally, 34,993 SNPs were retained for further analysis.

Experimental design and drought treatment

All experiments were conducted at the Plant Accelerator, Australian Plant Phenomics Facility, at the University of Adelaide, SA, Australia. The panel was phenotyped for a digital metric representing shoot growth over 20 days of progressive drought using an image-based phenomics platform. Each plant was phenotyped daily over a period of 20 days and the imaging interval was ~ 24 hr. The details of the experimental design are provided in Campbell *et al.* (2018). Briefly, each experiment consisted of 357 accessions from RDP1 and was repeated three times from February to April 2016. Two smart-greenhouses were used for each experiment. In each smart-greenhouse, the accessions were distributed across 432 pots positioned across 24 lanes. The experimental design used in this study is similar to the partially replicated paired design (Cullis *et al.* 2006). Here, a set of test entries is replicated while the remaining test entries are unreplicated. This design was slightly modified to accommodate the different water treatments and allow comparison of treatments within each accession. To this end, each accession was assigned to two consecutive pots, and the water treatments were randomly assigned to each pot. In each experiment, 54 accessions were randomly selected and replicated twice.

Seven days after transplant (DAT), plants were thinned to one seedling per pot. Two layers of blue mesh were placed on top of the pots to reduce evaporation. The plants were loaded on to the imaging system and were watered to 90% field capacity (FC) DAT. On the 13 DAT, each pot was watered to 90% and was imaged to obtain an initial phenotype before the onset of drought. One plant from each pair was randomly selected for drought treatment. Water was withheld from drought plants until 10% FC, and after which water was applied to maintain 10% FC. For the duration of the experiment, the control plants were maintained at 100% FC.

Statistical analysis of phenotypic data

Visible images were processed, and digital features were extracted using the open-source Python library Image Harvest (Knecht *et al.* 2016). The sum of plant pixels from the two sides and one top view of red/green/blue (RGB) images was summed and used as a measure of shoot biomass. This digital phenotype is referred to as the projected shoot area (PSA) throughout this study. Several studies have reported a high correlation between PSA estimates and shoot biomass (Campbell *et al.* 2015; Golzarian *et al.* 2011; Knecht *et al.* 2016). Prior to downstream analyses, outlier plants at each time point were detected for each trait using the 1.5 interquartile range rule, and potential outliers were plotted along with their treatment counterparts and inspected visually. Plants that exhibited abnormal growth patterns were removed. In total, we used records from 2,415 plants across both treatments. Briefly, there were 1,208-1,209 records in control for each day of imaging and 1,205-1,206 records for each day in low water treatment. All accessions had at least two records for each treatment-day block.

Raw phenotypic measurements were adjusted for downstream genetic analyses prior to fitting RRM. Best linear unbiased estimators (BLUE) were computed for each accession by fitting experimental effect with three levels and replication within experiment for some of the accessions. We postulated that observations at each time point follow the model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon}$, where \mathbf{X} and \mathbf{Z} are $n' \times f$ and $n' \times n$ orders of incident matrices linking observations (n') to systematic effects (f) and number of accessions (n), respectively, \mathbf{y} is an $n' \times 1$ vector of observations at each time point, $\boldsymbol{\beta}$ is a $f \times 1$ vector of systematic effects, \mathbf{u} is a $n \times 1$ vector of accession effects, and $\boldsymbol{\epsilon}$ is an $n' \times 1$ vector of residuals with $Var(\boldsymbol{\epsilon}) = \mathbf{I}\sigma_{\epsilon}^2$, where \mathbf{I} is an identity matrix. This was

followed by fitting a RRM-based genomic prediction approach to predict phenotypes as described below.

Random regression model

We conducted quantitative genetics modeling of image-derived phenotypes using a RRM to assess how well we could predict dynamic genetic signals. The RRM assumes that genetic effects and genetic variances are not constant and can vary continuously across the trajectory. This leads to better prediction of time-dependent complex traits by estimating heterogeneous single nucleotide polymorphism (SNP) effects across the trajectory. Specifically, we viewed the trajectory of digital image-processed longitudinal records as an infinite-dimensional characteristic that could be modeled by a smooth function (Meyer and Hill 1997; Van der Werf *et al.* 1998). Changes in PSA (BLUES) over time were modeled through Legendre polynomials and B-splines of time at phenotyping. The general formula for the RRM was as follows:

$$PSA_{ij} = \sum_k^{K_1} \phi(t)_{jk} \beta_k + \sum_k^{K_2} \phi(t)_{jk} u_{jk} + \sum_k^{K_3} \phi(t)_{jk} p_{jk} + \epsilon_{ij},$$

where $\phi(t)_{jk}$ is a time covariate coefficient defined by a basis function evaluated at time point t belonging to the j th accession; β_k is a k th fixed regression coefficient for the population's mean growth trajectory; u_{jk} is a k th random regression coefficient associated with the additive genetic effects of the j th accession; K_1 is the number of random regression parameters for fixed effect time trajectories; K_2 and K_3 are the number of random regression parameters for random effects; and p_{jk} is a k th permanent environmental random regression coefficient for the accession j . The starting values of index k , and K_1 , K_2 , and K_3 are defined separately for Legendre polynomials and B-splines below.

In the matrix notation, the above equation can be rewritten as

$$y = X\beta + Zu + Qpe + \epsilon,$$

where β is a vector of solutions for fixed regressions; u is the additive genetic random regression coefficients; pe is the permanent environmental random regression coefficients; ϵ is the residuals; and X , Z , and Q are corresponding incident matrices. Here, pe was defined as the resemblance between records of an individual due to non-random environmental effects that are persistent across the 20 time points. We assumed a multivariate-Gaussian distribution and the variance-covariance structure of

$$\text{Var} \begin{pmatrix} u \\ pe \\ \epsilon \end{pmatrix} = \begin{pmatrix} G \otimes C_u & 0 & 0 \\ 0 & I \otimes C_{pe} & 0 \\ 0 & 0 & R \end{pmatrix},$$

where $G = W_{sc}W'_{sc}/p$ is the genomic relationship matrix of VanRaden (2008), where W_{sc} represents a centered and standardized marker matrix and p is the number of markers; C_u is the covariance function between the random regression coefficients for the additive genetic effect; \otimes is the Kronecker product; C_{pe} is the covariance function between the random regression coefficients for the permanent environmental effects; and $R = I_n \sigma_{\epsilon(t)}^2$ is a diagonal matrix of heterogeneous residuals varying across times, where $\sigma_{\epsilon(t)}^2$ is the residual variance.

Choice of basis function

The choice of the basis function to model the shape of the longitudinal measurements is critical. An ideal basis function has adequate potential

to capture real patterns of changes in variance along a continuous scale (time) for a given trait (Meyer and Kirkpatrick 2005). In this study, we used RRM with two basis functions, *i.e.*, Legendre polynomials (Meyer 1998) and B-splines (Meyer 2005), to describe line-specific curves for the PSA trajectory over the day of imaging.

Legendre polynomials: Applying parametric shape functions for covariates of time is challenging because these covariates tend to generate high correlations among trajectories (Mrode 2014). For this reason, fitting Legendre polynomials of time at recording as covariables is a common choice to model growth curves because these polynomials greatly reduce the correlations between estimated random regression coefficients and make no prior assumptions regarding the shape of the longitudinal curve. This function has been used widely in animal breeding for many years (*e.g.*, Jamrozik and Schaeffer 1997) and has recently been used in plant breeding as well (Sun *et al.* 2017; Campbell *et al.* 2018; Marchal *et al.* 2019). Suppose d is the order of fit or degree of the polynomial. Legendre polynomials evaluated at the standardized time points were computed as $\Phi = M\Lambda$, where M is the t_{max} by $d + 1$ matrix containing the polynomials of the standardized time covariate $M_{k+1} = \left(\frac{2(t-t_{min})}{t_{max}-t_{min}} \right)^k - 1$

and Λ is the $d + 1 \times d + 1$ matrix of Legendre polynomial coefficients (Kirkpatrick *et al.* 1990). Here, $t_{min} = 1$ and $t_{max} = 20$ because PSA was measured for 20 days. We chose the same orders of polynomials for fixed, additive, and permanent environmental coefficients as previously described Schaeffer (2016). We compared linear ($k = 0, \dots, K_1 = K_2 = K_3 = 1$) and quadratic ($k = 0, \dots, K_1 = K_2 = K_3 = 2$) Legendre polynomials in this study. Thus, the numbers of regression coefficients were $d + 1 = 2$ and $d + 1 = 3$ for linear and quadratic Legendre polynomials, respectively.

B-splines: Spline functions consist of individual segments of polynomials joined at specific points called knots. B-splines first require determination of the total number of knots K . Although a large number of knots will increase complexity, too few knots will decrease accuracy. This basis function is reported to offer several advantages, including better numerical properties compared with polynomials, especially when there are high genetic variances at the extremes of the trajectory period, negative correlations between the most distant time point measurements, and a small number of records, particularly at the last stage of the trajectory (Meyer 2005; Misztal 2006). Here, we used equidistant knots, and the B-spline function was computed from Cox-de Boor's recursion formula (De Boor 2001). Given a preconsidered knot sequence of time t , the covariables for B-splines of degree $d = 0$ were defined by assuming values of unity for all points in a given interval or zero otherwise. For the i th interval given by knots

$$B_{i,d=0}(t) = \begin{cases} 1 & \text{if } T_i \leq t \leq T_{i+1} \\ 0 & \text{otherwise.} \end{cases}$$

where T is the threshold in time interval. According to De Boor (2001), the matrix Φ of B-spline for higher-order polynomials can be defined by recursion

$$B_{i,d}(t) = \frac{t - T_i}{T_{i+d} - T_i} B_{i,d-1}(t) + \frac{T_{i+d+1} - t}{T_{i+d+1} - T_{i+1}} B_{i+1,d-1}(t).$$

This indicates that a B-spline of degree d is simply a function of B-splines of degree $d - 1$. Note that the number of random regression coefficients depends on the number of knots and order of polynomials

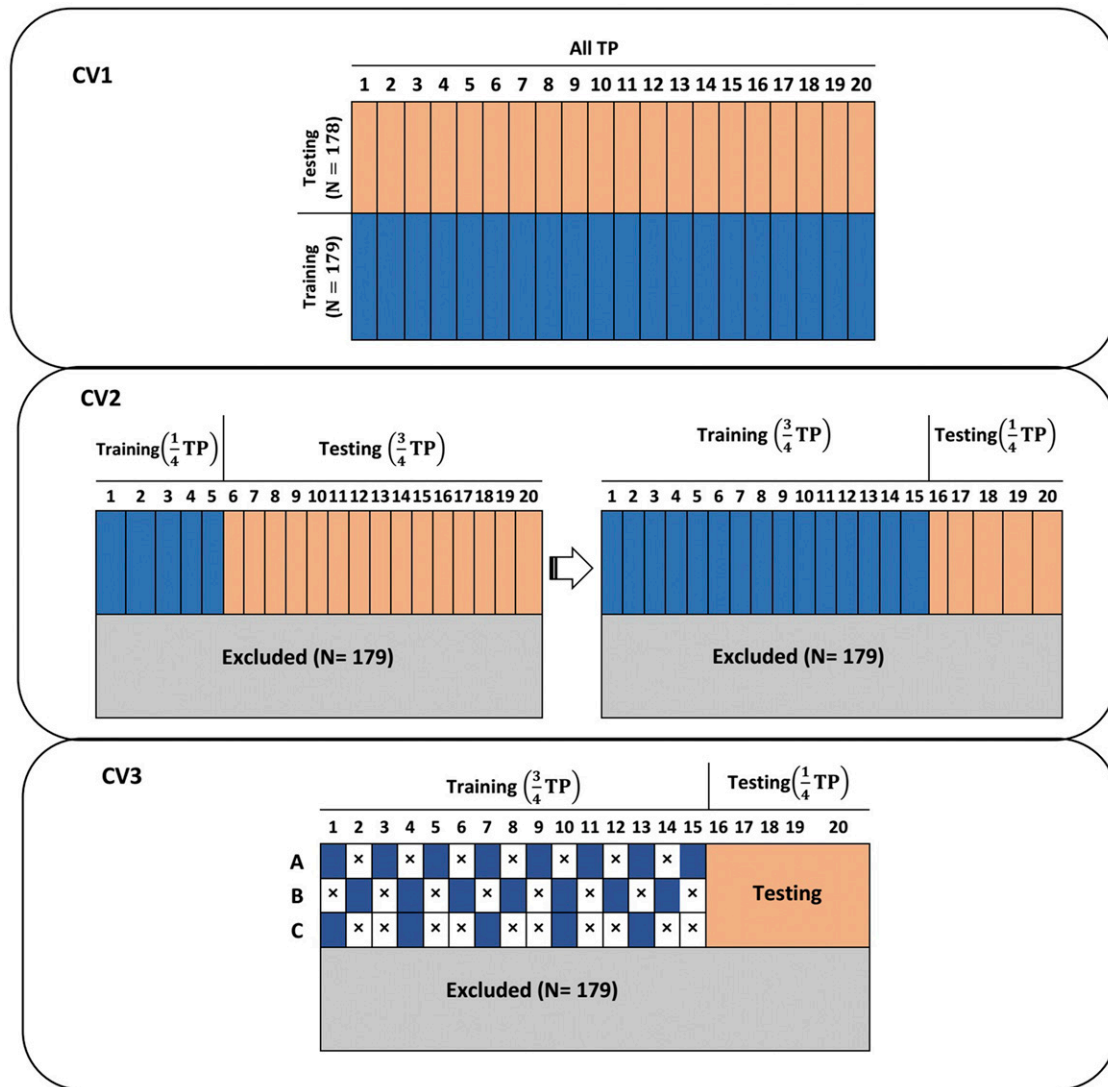


Figure 1 Pictorial representation of three cross-validation schemes used for predicting longitudinal projected shoot area (PSA) using a random regression model coupled with Legendre polynomials and B-splines. The data set consisted of 357 lines. CV1: 179 lines were used as the training set to predict PSA for the remaining 178 lines. Here, all 20 time points in the training set were used to predict PSA at each of 20 time points for a new set of lines. CV2: The data set was split into two longitudinal stages. The model was trained using the earlier growth stages to predict PSA at the second part of growth stages. We increased the number of time points used for training in a sequential manner. CV3: This was used to evaluate the impact of phenotyping frequency in the training data set on longitudinal prediction accuracy. Observations on odd days were used (CV3A), Observations on even days were used (CV3B), and keep one and delete two consecutive time points (CV3C). TP: time points.

for B-splines. In general, the number of regression coefficients is given by $K = s + d - 1$ (Meyer 2005). In this study, we fitted linear B-splines with $s = 3$ or $s = 4$ knots to divide the time points into equally spaced intervals. The same number of knots was considered for fixed trajectories, additive genetic, and permanent environmental coefficients. Thus, the numbers of regression coefficients were three ($k = 1, \dots, K_1 = K_2 = K_3 = 3 + 1 - 1 = 3$) and four ($k = 1, \dots, K_1 = K_2 = K_3 = 4 + 1 - 1 = 4$) for $s = 3$ and $s = 4$ knots, respectively.

Goodness of model fit

The goodness of fit of RRM was assessed by computing the Akaike's information criterion (AIC) (Akaike 1974) and the Schwarz–Bayesian information criterion (BIC) (Schwarz 1978). The best model was selected based on the largest AIC and BIC values after multiplying

by $-1/2$. We used the Wombat software to fit RRM in this study (Meyer 2007).

Cross-validation scenarios

As graphically represented in Figure 1, three different CV scenarios were designed to train the RRM. In all scenarios, prediction accuracy was evaluated by computing Pearson correlations between predicted genetic values and PSA in the testing set. Each of the CV scenarios is described below.

CV1: In the first CV scenario (CV1), the whole data set was divided into two subsets, *i.e.*, training and testing sets, each including 179 and 178 accessions, respectively. All 20 time points in the training set were fit to the RRM using Legendre polynomials and B-splines, and we predicted phenotypic values of 20 time points for lines in the testing set. Random

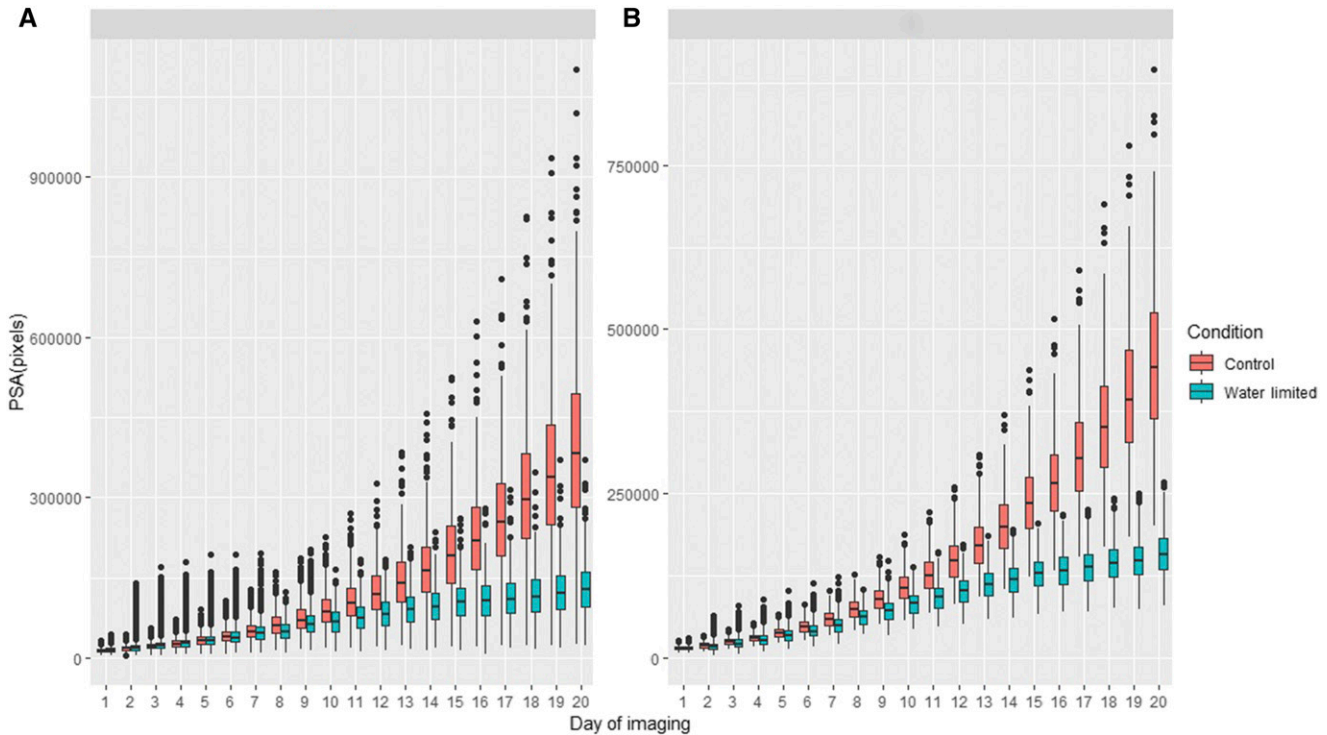


Figure 2 A: Box plots of projected shoot area (PSA) over the 20 days of imaging in two environments: controlled and water-limited conditions. B: Best linear unbiased estimators over the 20 days of imaging in two environments: controlled and water-limited conditions.

assignment of individuals into the training and testing sets was repeated 10 times. The equation for CV1 was set up in the following manner. The time-specific genetic value of the i th individual in the training set was computed as $\hat{\mathbf{g}}_{\text{tm}, i}^t = \Phi_t \mathbf{u}_i$, where $\hat{\mathbf{g}}_{\text{tm}, i}^t$ is the estimated genetic value of the individual i at time t ; Φ_t is the t th row vector of the $t_{\text{max}} \times K_1$ matrix Φ ; and \mathbf{u}_i is the i th column vector of the $K_1 \times n$ matrix \mathbf{u} . Then, a vector of predicted genetic values of individuals in the testing set at time t was obtained as $\hat{\mathbf{g}}_{\text{st}}^t = \mathbf{G}_{\text{st}, \text{tm}} \mathbf{G}_{\text{tm}, \text{tm}}^{-1} \hat{\mathbf{g}}_{\text{tm}}^t$, where $\mathbf{G}_{\text{st}, \text{tm}}$ is the genomic relationship matrix between the testing and training set and $\mathbf{G}_{\text{tm}, \text{tm}}^{-1}$ is the inverse of genomic relationship matrix between the training set. Because CV1 is not a forecasting task, a standard multi-trait model (MTM) was also fitted as a baseline model considering longitudinal traits as different traits (Henderson and Quaas 1976).

$$\mathbf{y} = \mathbf{X}\beta + \mathbf{Z}\mathbf{u} + \epsilon,$$

where β is a vector of systematic effects; \mathbf{u} is the vector of additive genetic values; ϵ is the residuals; and \mathbf{X} and \mathbf{Z} are corresponding incident matrices. The joint distribution of \mathbf{u} and ϵ follows multivariate normal

$$\begin{pmatrix} \mathbf{u} \\ \epsilon \end{pmatrix} \sim N \left[\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \begin{pmatrix} \mathbf{G} \otimes \Sigma_u & 0 \\ 0 & \mathbf{I} \otimes \Sigma_\epsilon \end{pmatrix} \right],$$

where \mathbf{G} is the genomic relationship matrix, \mathbf{I} is the identity matrix, and Σ_u and Σ_ϵ refer to 20×20 dimensional additive genetic and residual variance-covariance matrices, respectively. In other words, the diagonal and off-diagonal elements of Σ include variance at each time point and covariance between time points, respectively. The BLUPF90 family of programs was used to fit MTM with 20 traits (Misztal *et al.* 2002).

CV2: The second CV scenario (CV2) was related to forecasting future phenotypes from longitudinal traits at early time points. Individuals in

the training set were used to forecast their yet-to-be observed PSA values at later time points from information available at earlier time points. The first quarter of the time points $\{t = 1, 2, 3, 4, 5\}$ was used as the training set, and the remaining time points $\{t = 6, 7, \dots, 20\}$ were predicted for each line in the training set. This was followed by sequentially increasing the number of time points used to train the model so that in the last run, three quarters of the time points $\{t = 1, 2, \dots, 15\}$ were used in the training set to forecast phenotypes at the last quarter of time points $\{t = 16, 17, 18, 19, 20\}$. This CV scenario was designed to find a sufficient set of earlier time points to obtain reasonable longitudinal prediction accuracy and is known as walk forward validation. We set up the CV2 equation by first estimating the random regression coefficient matrix \mathbf{u} using $\Phi_{1:t}$, which was computed from time point 1 to time point t . The prediction of future time points t' ($t + 1 \leq t' \leq t_{\text{max}}$) for an individual i was carried out by $\hat{\mathbf{g}}^{t'} = \Phi_{t'} \mathbf{u}_i$, where $\Phi_{t'}$ is the t' th row vector of $t_{\text{max}} - t$ by $K + 1$ matrix Φ ; and \mathbf{u}_i is the i th column vector of the number of random regression coefficients by n matrix \mathbf{u} .

CV3: The third CV scenario (CV3) was designed to evaluate whether it was possible to reduce the phenotyping frequency while still maintaining a high prediction accuracy for the last quarter of observations. We used the last case in CV2 such that time points $\{t = 1, 2, \dots, 15\}$ were used for the training set to forecast the last quarter of observations $\{t = 16, 17, 18, 19, 20\}$. We then reduced the number of time points used in the training set as follows: A, observations on odd days $\{t = 1, 3, \dots, 15\}$ were used; B, observations on even days $\{t = 2, 4, \dots, 14\}$ were used; C, keep one and delete two consecutive time points. In CV2 and CV3 scenarios, half of the individuals were randomly selected to fit and evaluate the model, and the analysis was repeated 10 times. If the loss of prediction accuracy was minimal, it was possible to reduce the phenotyping cost. The equation for CV3 was set up in the same way as that for CV2.

Table 1 Assessing goodness of fit for two random regression models (Legendre polynomials and B-splines) used to predict projected shoot area measured over 20 time points

Condition	BF	Log L	-0.5 AIC	-0.5 BIC	<i>p</i>
CON	LEGL	-32414.493	-32440.493	-32529.839	26
	LEGO	-32412.550	-32444.550	-32554.512	32
	BSPL3	-32408.862	-32440.862	-32550.824	32
	BSPL4	-32404.142	-32444.142	-32581.592	40
WL	LEGL	-26011.867	-26037.867	-26127.213	26
	LEGO	-26009.267	-26041.267	-26151.229	32
	BSPL3	-26006.205	-26038.205	-26148.167	32
	BSPL4	-26005.537	-26045.537	-26182.986	40

CON: control environment; WL: water-limited environment; BF: basis function; LEGL: Legendre polynomial linear; LEGO: Legendre polynomial quadratic; BSPL3: B-spline linear with three knots; BSPL4: B-spline linear with four knots; Log L: log likelihood; AIC: Akaike information criterion; BIC: Bayesian information criterion; and *p*: number of parameters.

Data availability

Phenotypic data used herein are available in Supplementary File S1 at Figshare. Genotypic data regarding the rice accessions can be downloaded from the rice diversity panel website (<http://www.ricediversity.org/>) and also available in Supplementary File S2 at Figshare. Supplemental material available at FigShare: <https://doi.org/10.25387/g3.9383543>.

RESULTS

Assessing model fit

Figures 2A and 2B show the box plots of the original PSA and BLUE for the phenotypic trajectories over the 20 days of imaging for control and water-limited conditions. The PSA for control and water-limited plants

diverged significantly six days after initiating the drought treatment ($p < 0.0025$). Supplementary Figure 1 (File S3) shows the linear or quadratic forms of Legendre polynomials and three and four knot-based B-spline curves over 20 days of imaging. For Legendre polynomials, intercept, linear, and quadratic coefficients are represented in black, red, and green, respectively. For B-spline, knot 1, knot 2, and knot 3 are represented in black, red, and green, respectively.

Table 1 summarizes the goodness of fits of RRM coupled with linear and quadratic Legendre polynomials and B-spline functions in control and water-limited conditions. For the Legendre polynomials, quadratic forms require more parameters to be estimated compared with linear forms. Similar to observation for B-splines, the presence of a greater number of knots suggested that there were more parameters to be estimated. Under control conditions, the best goodness of fit was obtained by linear Legendre polynomials, followed by linear B-splines with three knots, linear B-splines with four knots, and quadratic Legendre polynomials according to AIC scores. According to BIC scores, linear Legendre polynomials performed the best, followed by linear B-splines with three knots, quadratic Legendre polynomials, and linear B-splines with four knots. Under water-limited conditions, the best goodness of fit was given by linear Legendre polynomials, followed by linear B-splines with three knots, quadratic Legendre polynomials, and linear B-splines with four knots for both AIC and BIC scores. The number of parameters in the model varied from 26 to 40.

Cross-validation

The results from CV1 are shown in Figure 3. This CV was designed to evaluate the accuracy of predicting testing set individuals using all time points. Under control conditions, MTM performed relatively

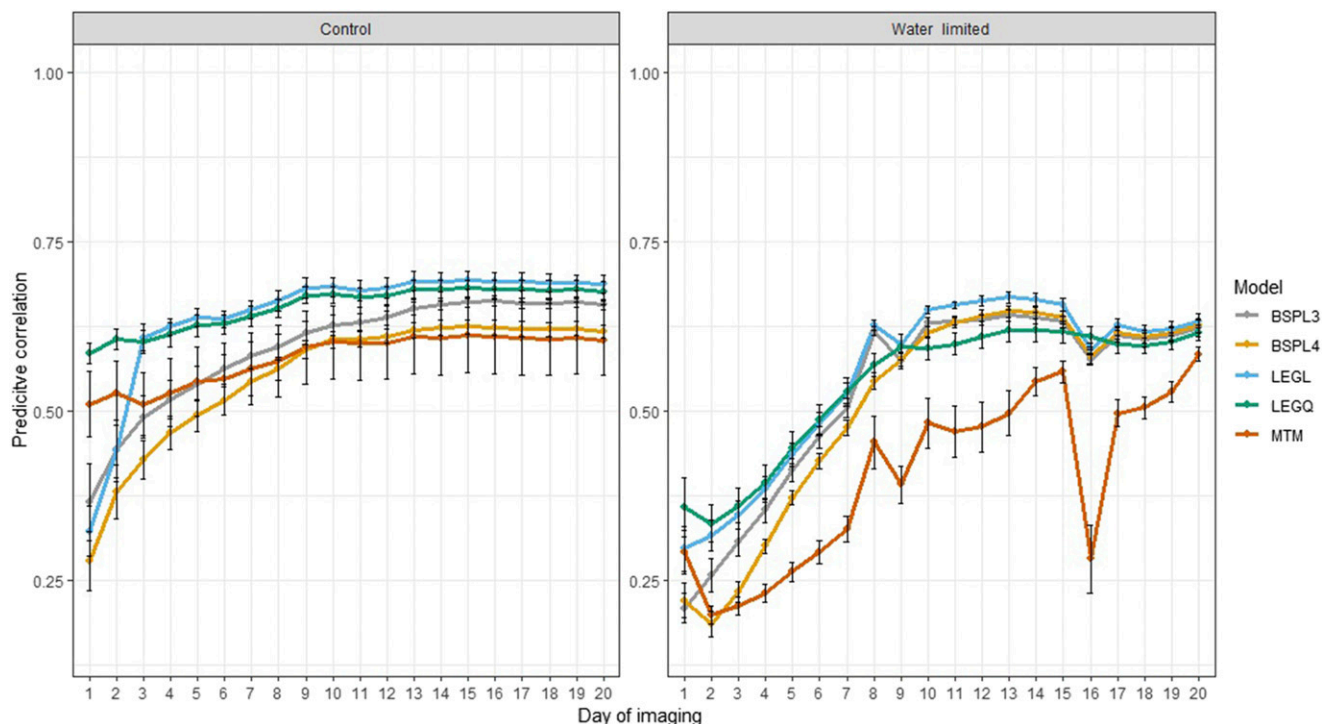


Figure 3 Prediction accuracy obtained from cross-validation 1 scenario. Total of 179 lines were used as the training set to predict PSA for the remaining 178 lines. Here, all 20 time points in the training set were used to predict PSA at each of 20 time points for a new set of lines. LEGL: linear Legendre polynomials; LEGQ: quadratic Legendre polynomials; BSPL3: linear B-splines with three knots; BSPL4: linear B-spline with four knots; MTM: multi-trait model.

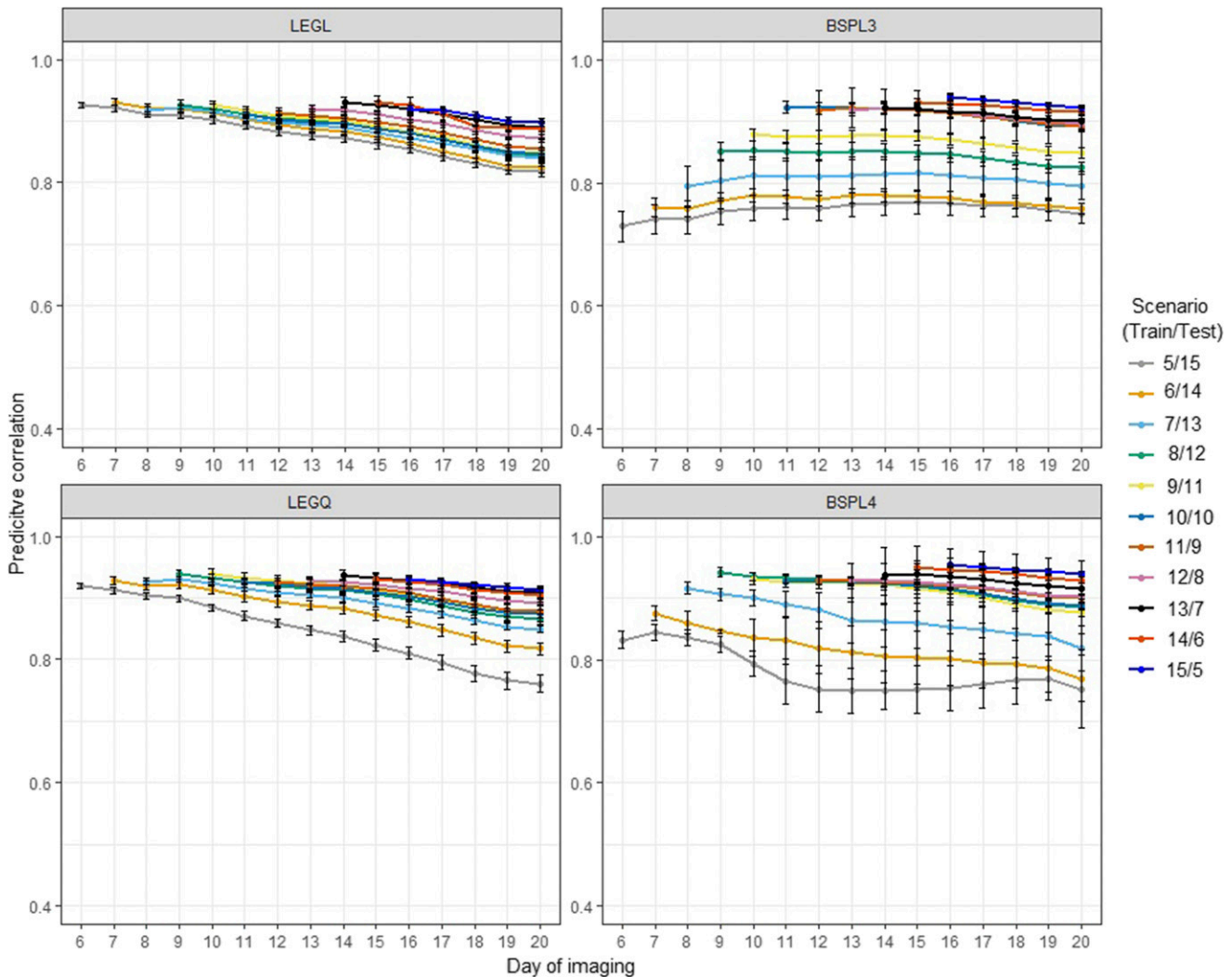


Figure 4 Prediction accuracy of cross-validation scenario 2 in control conditions. Each line depicts the different number of training and testing sets partitioning at the time point levels. LEGL: linear Legendre polynomials; LEGQ: quadratic Legendre polynomials; BSPL3: linear B-splines with three knots; BSPL4: linear B-spline with four knots.

better than RRM up to day 3. The prediction accuracy of RRM increased subsequently and after the 10th day of imaging, the best prediction was given by linear Legendre, followed by quadratic Legendre, linear B-spline with three knots, and linear B-spline with four knots. Overall, RRM performed better than MTM, and linear Legendre was the best prediction machine throughout the growth stages. Under water-limited conditions, prediction accuracy was lower compared with those of control conditions. All RRM delivered higher prediction than MTM except for the first two time points. Although Legendre polynomials performed better than B-splines until day 7, the difference between these approaches became negligible afterward.

Figures 4 and 5 show the accuracy of CV2 under control and water-limited conditions, respectively. This CV was designed to test how much information from the previous time points was required to achieve reasonable prediction accuracy at later growth stages. Under control conditions, we found that the best prediction for the last five time points was achieved when using information from all prior time points (15/5 CV2 subscenario). This suggested that having more information from previous time points to train the model would result in

higher prediction accuracy. Using the first five time points to train the model resulted in the worse prediction (5/15 CV2 subscenario). Thus, it is likely that the prediction accuracy in RRM declined because we attempted to estimate numerous parameters from only five time points. Legendre polynomials yielded better and more stable prediction than B-splines. We observed a similar trend under water-limited conditions; that is, using more previous time points to train the model resulted in higher prediction accuracy. However, the accuracy of prediction was unstable and decreased dramatically. There was no noticeable difference between the Legendre polynomials and B-splines in terms of performance.

Figures 6 and 7 show the CV3 accuracy under control and water-limited conditions, respectively. We designed this CV to evaluate whether it was possible to reduce phenotyping frequency and phenotyping costs without sacrificing prediction accuracy. Under control conditions, the prediction accuracy of CV3A, CV3B, and CV3C all decreased relative to the benchmark scenario in CV2, where all of the first 15 time points were used for the training set to forecast the last five time points. Although removing two consecutive time points did not improve performance (CV3C), the prediction accuracy from

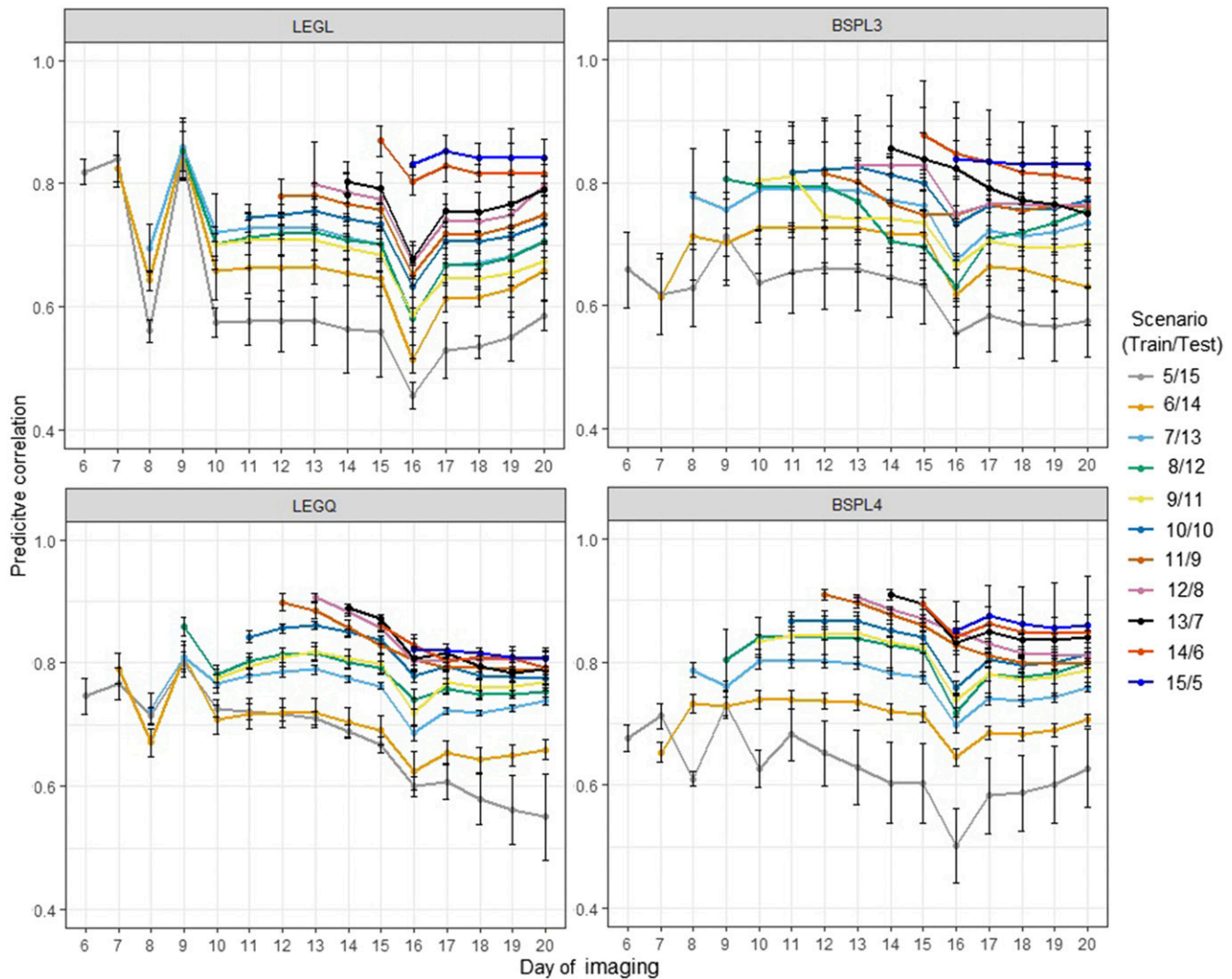


Figure 5 Prediction accuracy of cross-validation scenario 2 in water-limited conditions. Each line depicts the different number of training and testing sets partitioning at the time point levels. LEGL: linear Legendre polynomials; LEQG: quadratic Legendre polynomials; BSPL3: linear B-splines with three knots; BSPL4: linear B-spline with four knots.

phenotyping every other day was still relatively high (CV3A and CV3B). In general, the linear B-splines performed the best, and differences between scenarios were minimal. Under water-limited conditions, we observed the same trend, but the prediction accuracy was more unstable and decreased relative to control conditions. The quadratic Legendre polynomials and B-splines with four knots did not perform well, possibly due to overfitting.

DISCUSSION

Image-based automated HTP technologies offer great potential for characterizing multi-faceted phenotypes at high temporal resolution. The use of HTP platforms plays a pivotal role in accelerating breeding efforts by providing the temporal resolution needed for capturing adaptive responses to environmental challenges, but the development of statistical methodologies to analyze image-based function-valued phenotypes has not kept pace with our ability to generate HTP data. Because phenomics and genomics landscapes for plants are constantly advancing, parallel efforts are required to develop tools for integrating diverse genomic and phenomic datasets characterized by high temporal resolution in genetic analysis. Rice is one of the most

drought sensitive cereal crops, resulting in substantial yield losses. With predictions for greater climatic shifts in the future and increasing competition for fresh water resources, research that leverages the full potential of genomics and phenomics is needed to elucidate the genetic and physiological basis of drought tolerance. However, there is currently a lack of information regarding the modeling of temporal HTP data.

RRM identifies the effects of heterogeneous SNPs that transiently influence key traits and translates this to prediction of phenotypes. The main idea behind RRM is to describe subject-specific curves through basis functions (Meyer and Kirkpatrick 2005). Although RRM has been successfully applied to pedigree-based animal breeding (Schaeffer and Jamrozik 2008), its utility is largely limited to evaluating goodness-of-fit for candidate models rather than CV-based prediction, and its integration into HTP data has not been reported. In this study, we coupled HTP data with high-density genomic information to carry out longitudinal prediction by capturing time-specific genetic signals. A diverse panel of rice accessions subjected to drought stress was used to illustrate the utility of the RRM for evaluating Legendre polynomials and B-splines of time at recording.

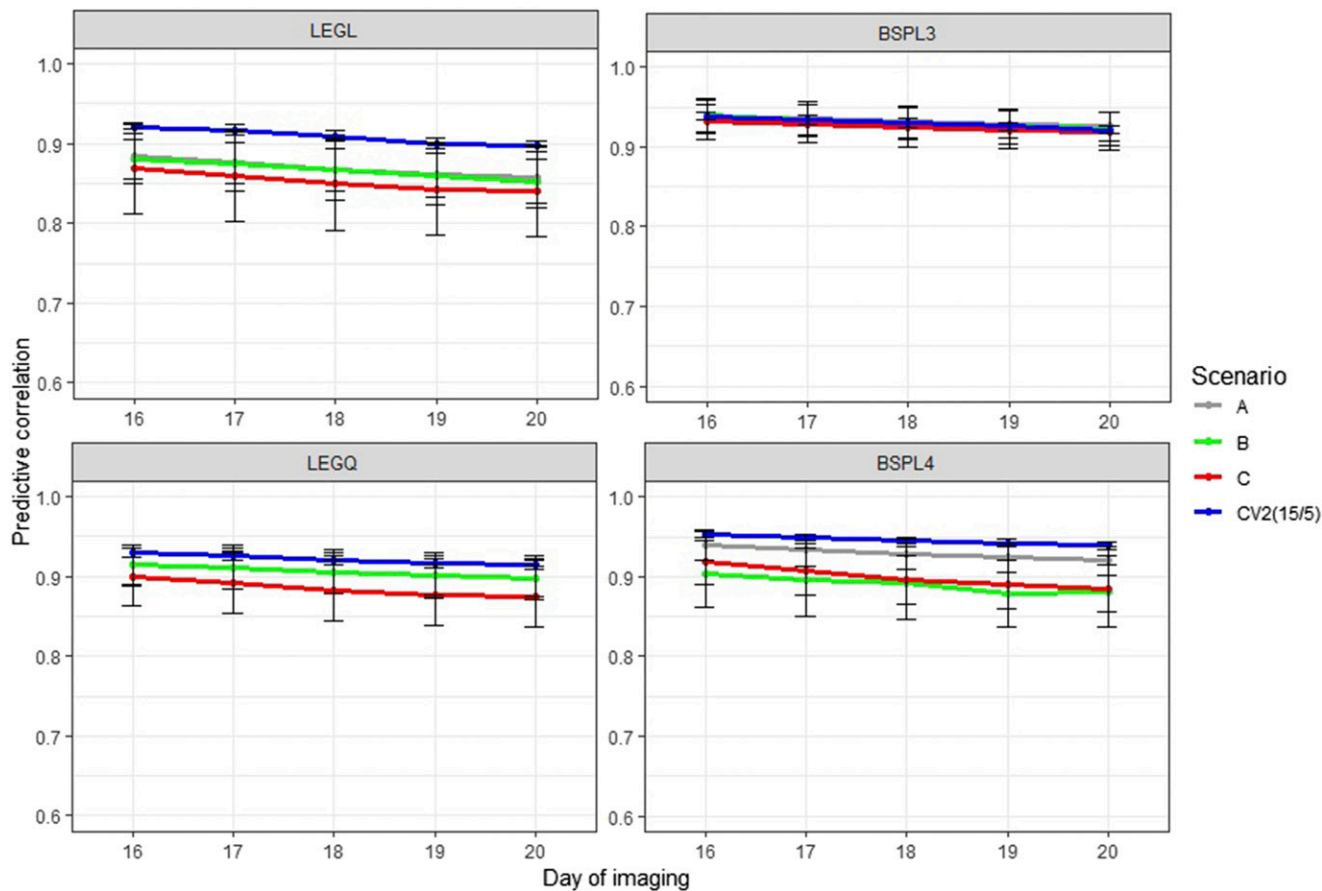


Figure 6 Prediction accuracy of cross-validation scenario 3 in control conditions. A: only observations in the odd days were used; B: only observations in the even days were used; C: keep one and delete two consecutive time points; CV2: use all available previous time points; LEGL: linear Legendre polynomials; LEGQ: quadratic Legendre polynomials; BSPL3: linear B-splines with three knots; BSPL4: linear B-spline with four knots.

Longitudinal prediction

We found that it was possible to model longitudinal PSA responses under water-limited conditions, albeit with decreased prediction accuracy compared with that of the control. We also placed particular emphasis on comparing two basis functions, *i.e.*, Legendre polynomials and B-splines. To the best of our knowledge, the current study is the first to use a B-spline function to evaluate longitudinal prediction accuracy in the RRM applied to HTP data. Linear B-spline functions with $s = 3$ (two segments) or $s = 4$ knots (three segments) were used. B-splines have been reported to have better numerical properties (*e.g.*, lower computational requirement and faster convergence) than Legendre polynomials because each coefficient of a function affects only a part of the trajectory and can be used to estimate genetic parameters more smoothly while still adequately capturing the features of dynamic data (Iwaisaki *et al.* 2005; Baldi *et al.* 2010).

We observed differences in prediction accuracy across models during early growth stages; however, differences were incremental when predicting later growth stages in the CV1 scenario, in which the training and testing sets were partitioned based on individuals. Overall, linear Legendre polynomials performed the best and was clearly an advancement over the MTM. Prediction performance in CV2, in which the training and testing sets were partitioned according to growth stages rather individuals, showed that it was possible to predict future phenotypes from information available from earlier trajectories. Here,

linear and quadratic Legendre polynomials produced the highest and most stable prediction accuracy under control conditions, whereas linear B-splines with three knots performed better in the water-limited environment. The final scenario (CV3) demonstrated that we could decrease the phenotyping frequency by only phenotyping every other day to reduce the phenotyping cost while minimizing the loss of prediction accuracy. In this case, linear B-spline with three knots performed relatively well.

B-spline functions require two parameters (the position of the knots and the number of knots) to be tuned. The position of knots can be chosen based on a trajectory pattern such that more knots are placed for rapidly changing time points, whereas less knots are placed for time points with slow changes (Misztal 2006). Thus, the position of knots should be carefully chosen if the number of phenotyped individuals varies substantially at each growth stage. We chose equidistant knots in the current study because all accessions were phenotyped on the same days during the trajectory. The number of knots determines the number of segments fitted. When more knots are specified, the model becomes more complex. Although we used $s = 3$ and $s = 4$ based on previous literature and a visual inspection of the observed phenotypic trajectory, further investigations are warranted to explore the impact of the number of knots on prediction accuracy. The performance of quadratic B-spline functions was not evaluated in the current study because we encountered convergence issues, possibly due to the small sample size. In general, we found that the advantages of B-splines in inferential

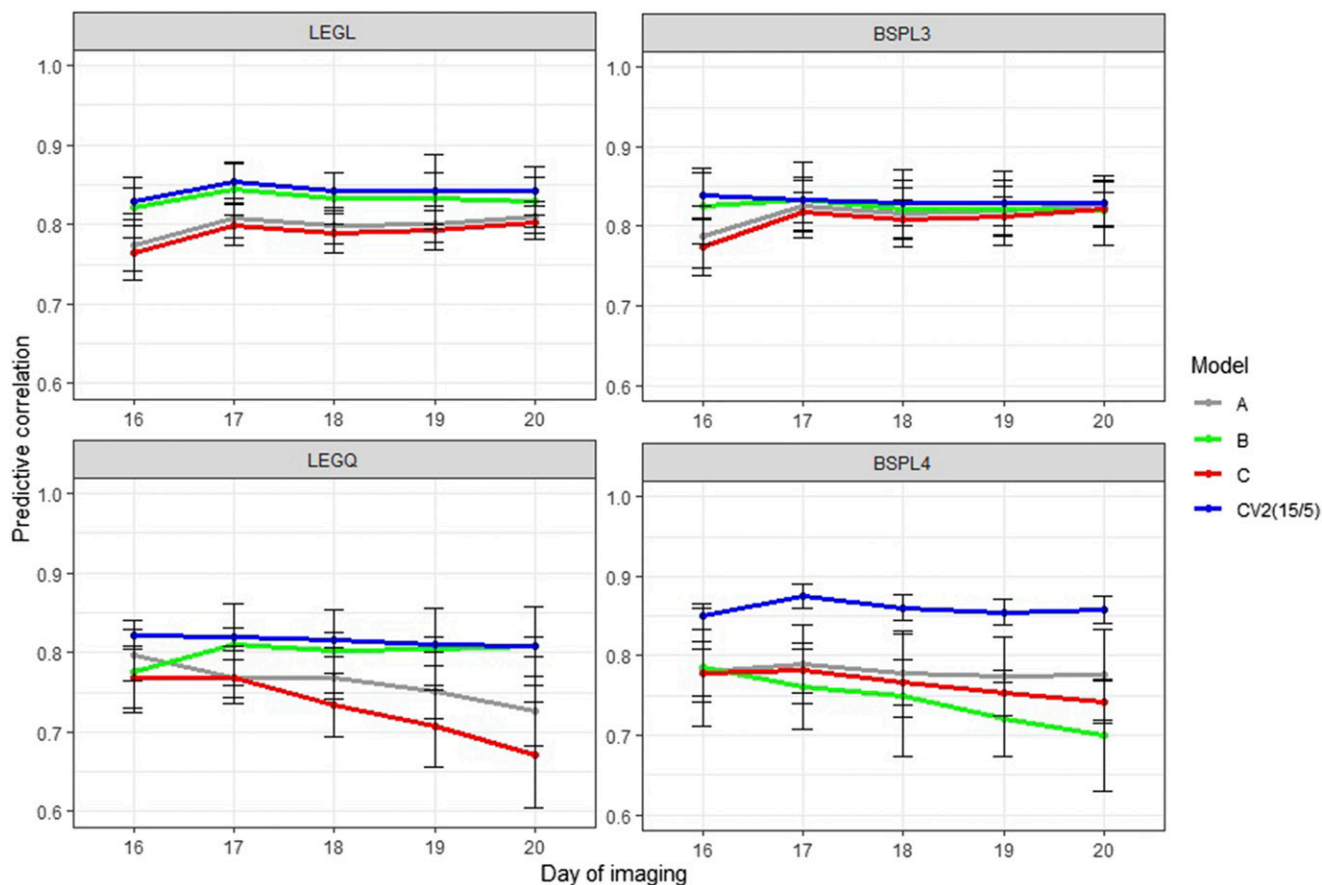


Figure 7 Prediction accuracy of cross-validation scenario 3 in water-limited conditions. A: only observations in the odd days were used; B: only observations in the even days were used; C: keep one and delete two consecutive time points; CV2: use all available previous time points; LEGL: linear Legendre polynomials; LEGQ: quadratic Legendre polynomials; BSPL3: linear B-splines with three knots; BSPL4: linear B-spline with four knots.

tasks compared with Legendre polynomials were not shown clearly in terms of prediction. This is likely because PSA trajectories were relatively simple exponential or monotonically increasing trajectories without obvious inflection points, indicating that the potential of B-splines was not able to be fully exploited in the current study.

Choice of parameters

We also found that ranking the models according to AIC and BIC revealed only mild agreement with prediction performance evaluated by CV, suggesting that the RRM that gives the best goodness-of-fit is not guaranteed to deliver the best prediction and vice versa. The choice for the order of fit or the number of knots is arguably the most challenging modeling aspect in the RRM. In the majority of literature describing the RRM, this parameter is mainly chosen based on AIC, BIC, or the eigendecomposition of the covariance matrix. The major issue regarding this approach is that there is a tendency to simply pick a model with the highest order of fit or the largest number of knots. However, this study, suggests finding the best parameter in terms of prediction accuracy obtained from CV.

Future perspective

We anticipate that the current work will guide us to conduct genomic selection of economically important traits on the longitudinal scale for the purpose of breeding crops that are better adapted to new environments or to less favorable challenging climatic conditions. Although in

the current study, our aim was to assess RRM for genomic prediction of shoot biomass under contrasting water regimes, these frameworks can be extended to any time-resolved phenotype, provided there are enough time points with complete or partial records. Owing to the accessibility of HTP platforms in the public sector as well as the growing availability of unmanned aerial vehicles and other autonomous field-based platforms, many breeding programs are currently generating temporal phenotypic data. Although the temporal phenotypes themselves may not be a target of selection, these data can be utilized to improve selection for conventional end-point phenotypes such as yield. For instance, Sun *et al.* (2017) used parameters from RRM as covariates in a mixed model to improve prediction for yield in drought-stressed environments in wheat.

The identification of genomic components over trajectories will provide information regarding the optimum time points to maximize cost-effective selection or to design a genome-assisted breeding program aiming to change the shape of the longitudinal response curve (Schaeffer 2004). Using our approach, we could evaluate all changes in plant biomass accumulation during the course of the experiment, in contrast to single time point analyses. Thus, we expect that RRM analysis will become the norm for modeling trajectories of function-valued HTP data because such approaches could be considered an extension of the widely used genomic best linear unbiased prediction model for time series data. Lastly, the utility of the RRM does not preclude its use in other applications.

For example, the RRM offers a new avenue for performing longitudinal GWAS (e.g., Howard *et al.* 2015; Campbell *et al.* 2019) and genotype-by-environment interactions using the reaction norm (Arnold *et al.* 2019). In summary, an RRM using Legendre polynomial or spline functions could be an effective option for modeling trait trajectories of HTP data and may have potential applications in characterizing phenotypic plasticity in plants.

ACKNOWLEDGMENTS

This work was supported by the National Science Foundation under Grant Number 1736192 to HW and GM, and Virginia Polytechnic Institute and State University startup funds to GM. MTC and HW designed and conducted the experiments. MM analyzed the data. MM and GM conceived the idea and wrote the manuscript. MTC and HW discussed results and revised the manuscript. GM supervised and directed the study. All authors read and approved the manuscript.

LITERATURE CITED

- Akaike, H., 1974 A new look at the statistical model identification, pp. 215–222 in *Selected Papers of Hirotugu Akaike*, Springer, Berlin. https://doi.org/10.1007/978-1-4612-1694-0_16
- Arnold, P. A., L. E. Kruuk, and A. B. Nicotra, 2019 How to analyse plant phenotypic plasticity in response to a changing climate. *New Phytol.* 222: 1235–1241. <https://doi.org/10.1111/nph.15656>
- Baldi, F., M. Alencar, and L. G. Albuquerque, 2010 Random regression analyses using b-splines functions to model growth from birth to adult age in canchim cattle. *J. Anim. Breed. Genet.* 127: 433–441. <https://doi.org/10.1111/j.1439-0388.2010.00873.x>
- Browning, S. R., and B. L. Browning, 2007 Rapid and accurate haplotype phasing and missing-data inference for whole-genome association studies by use of localized haplotype clustering. *Am. J. Hum. Genet.* 81: 1084–1097. <https://doi.org/10.1086/521987>
- Campbell, M. T., A. C. Knecht, B. Berger, C. J. Brien, D. Wang *et al.*, 2015 Integrating image-based phenomics and association analysis to dissect the genetic architecture of temporal salinity responses in rice. *Plant Physiol.* 168: 1476–1489. <https://doi.org/10.1104/pp.15.00450>
- Campbell, M. T., H. Walia, and G. Morota, 2018 Utilizing random regression models for genomic prediction of a longitudinal trait derived from high-throughput phenotyping. *Plant Direct* 2: 1–11. <https://doi.org/10.1002/pld3.80>
- Campbell, M., M. Momen, H. Walia, and G. Morota, 2019 Leveraging breeding values obtained from random regression models for genetic inference of longitudinal traits. *Plant Genome* 12: 180075. <https://doi.org/10.3835/plantgenome2018.10.0075>
- Cullis, B. R., A. B. Smith, and N. E. Coombes, 2006 On the design of early generation variety trials with correlated data. *J. Agric. Biol. Environ. Stat.* 11: 381–393. <https://doi.org/10.1198/108571106X154443>
- De Boor, C., 2001 *A Practical Guide to Splines*, Vol. 27, Revised Edition. Springer-Verlag, New York.
- Furbank, R. T., and M. Tester, 2011 Phenomics-technologies to relieve the phenotyping bottleneck. *Trends Plant Sci.* 16: 635–644. <https://doi.org/10.1016/j.tplants.2011.09.005>
- Ge, Y., G. Bai, V. Stoerger, and J. C. Schnable, 2016 Temporal dynamics of maize plant growth, water use, and leaf water content using automated high throughput rgb and hyperspectral imaging. *Comput. Electron. Agric.* 127: 625–632. <https://doi.org/10.1016/j.compag.2016.07.028>
- Golzarian, M. R., R. A. Frick, K. Rajendran, B. Berger, S. Roy *et al.*, 2011 Accurate inference of shoot biomass from high-throughput images of cereal plants. *Plant Methods* 7: 2. <https://doi.org/10.1186/1746-4811-7-2>
- Henderson, C., and R. Quaas, 1976 Multiple trait evaluation using relatives' records. *J. Anim. Sci.* 43: 1188–1197. <https://doi.org/10.2527/jas1976.4361188x>
- Howard, J. T., S. Jiao, F. Tiezzi, Y. Huang, K. A. Gray *et al.*, 2015 Genome-wide association study on legendre random regression coefficients for the growth and feed intake trajectory on Duroc Boars. *BMC Genomics* 16: 59. <https://doi.org/10.1186/s12863-015-0218-8>
- Iwaisaki, H., S. Tsuruta, I. Misztal, and J. Bertrand, 2005 Genetic parameters estimated with multitrait and linear spline-random regression models using gelbvieh early growth data. *J. Anim. Sci.* 83: 757–763. <https://doi.org/10.2527/2005.834757x>
- Jamrozik, J., and L. Schaeffer, 1997 Estimates of genetic parameters for a test day model with random regressions for yield traits of first lactation holsteins. *J. Dairy Sci.* 80: 762–770. [https://doi.org/10.3168/jds.S0022-0302\(97\)75996-4](https://doi.org/10.3168/jds.S0022-0302(97)75996-4)
- Kirkpatrick, M., D. Lofsvold, and M. Bulmer, 1990 Analysis of the inheritance, selection and evolution of growth trajectories. *Genetics* 124: 979–993.
- Knecht, A. C., M. T. Campbell, A. Caprez, D. R. Swanson, and H. Walia, 2016 Image Harvest: an open-source platform for high-throughput plant image processing and analysis. *J. Exp. Bot.* 67: 3587–3599. <https://doi.org/10.1093/jxb/erw176>
- Marchal, A., C. D. Schlichting, R. Gobin, P. Balandier, F. Millier *et al.*, 2019 Deciphering hybrid larch reaction norms using random regression. *G3: Genes, Genomes, Genetics* 9: 21–32.
- Meyer, K., 1998 Estimating covariance functions for longitudinal data using a random regression model. *Genet. Sel. Evol.* 30: 221. <https://doi.org/10.1186/1297-9686-30-3-221>
- Meyer, K., 2005 Random regression analyses using B-splines to model growth of australian angus cattle. *Genet. Sel. Evol.* 37: 473. <https://doi.org/10.1186/1297-9686-37-6-473>
- Meyer, K., 2007 Wombat - A tool for mixed model analyses in quantitative genetics by restricted maximum likelihood (reml). *J. Zhejiang Univ. Sci. B* 8: 815–821. <https://doi.org/10.1631/jzus.2007.B0815>
- Meyer, K., and W. B. Hill, 1997 Estimation of genetic and phenotypic covariance functions for longitudinal or repeated records by restricted maximum likelihood. *Livest. Prod. Sci.* 47: 185–200. [https://doi.org/10.1016/S0301-6226\(96\)01414-5](https://doi.org/10.1016/S0301-6226(96)01414-5)
- Meyer, K., and M. Kirkpatrick, 2005 Up hill, down dale: quantitative genetics of curvaceous traits. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* 360: 1443–1455. <https://doi.org/10.1098/rstb.2005.1681>
- Misztal, I., 2006 Properties of random regression models using linear splines. *J. Anim. Breed. Genet.* 123: 74–80. <https://doi.org/10.1111/j.1439-0388.2006.00582.x>
- Misztal, I., S. Tsuruta, T. Strabel, B. Auvray, T. Druet *et al.*, 2002 Blupf90 and related programs (bgf90). In *Proceedings of the 7th World Congress on Genetics Applied to Livestock Production*, volume 33. 743–744.
- Mrode, R. A., 2014 *Linear models for the prediction of animal breeding values*, CABI, Oxfordshire, UK. <https://doi.org/10.1079/9781780643915.0000>
- Purcell, S., B. Neale, K. Todd-Brown, L. Thomas, M. A. Ferreira *et al.*, 2007 Plink: a tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* 81: 559–575. <https://doi.org/10.1086/519795>
- Schaeffer, L., 2016 *Random regression models*. Available in <http://animal-biosciences.uoguelph.ca/~lrs/BOOKS/rrmbook.pdf>.
- Schaeffer, L., and J. Jamrozik, 2008 Random regression models: a longitudinal perspective. *J. Anim. Breed. Genet.* 125: 145–146. <https://doi.org/10.1111/j.1439-0388.2008.00748.x>
- Schaeffer, L. R., 2004 Application of random regression models in animal breeding. *Livest. Prod. Sci.* 86: 35–45. [https://doi.org/10.1016/S0301-6226\(03\)00151-9](https://doi.org/10.1016/S0301-6226(03)00151-9)
- Schwarz, G., 1978 Estimating the dimension of a model. *Ann. Stat.* 6: 461–464. <https://doi.org/10.1214/aos/1176344136>
- Sun, J., J. E. Rutkoski, J. A. Poland, J. Crossa, J.-L. Jannink *et al.*, 2017 Multitrait, random regression, or simple repeatability model in high-throughput phenotyping data improve genomic prediction for wheat grain yield. *Plant Genome* 10. <https://doi.org/10.3835/plantgenome2016.11.0111>

- Tester, M., and P. Langridge, 2010 Breeding technologies to increase crop production in a changing world. *Science* 327: 818–822. <https://doi.org/10.1126/science.1183700>
- Van der Werf, J. H., M. E. Goddard, and K. Meyer, 1998 The use of covariance functions and random regressions for genetic evaluation of milk production based on test day records. *J. Dairy Sci.* 81: 3300–3308. [https://doi.org/10.3168/jds.S0022-0302\(98\)75895-3](https://doi.org/10.3168/jds.S0022-0302(98)75895-3)
- VanRaden, P. M., 2008 Efficient methods to compute genomic predictions. *J. Dairy Sci.* 91: 4414–4423. <https://doi.org/10.3168/jds.2007-0980>
- White, I., R. Thompson, and S. Brotherstone, 1999 Genetic and environmental smoothing of lactation curves with cubic splines. *J. Dairy Sci.* 82: 632–638. [https://doi.org/10.3168/jds.S0022-0302\(99\)75277-X](https://doi.org/10.3168/jds.S0022-0302(99)75277-X)
- Zhao, K., C.-W. Tung, G. C. Eizenga, M. H. Wright, M. L. Ali *et al.*, 2011 Genome-wide association mapping reveals a rich genetic architecture of complex traits in *Oryza sativa*. *Nat. Commun.* 2: 467. <https://doi.org/10.1038/ncomms1467>

Communicating editor: J. Holland