

1 **Predicting longitudinal traits derived from**
2 **high-throughput phenomics in contrasting**
3 **environments using genomic Legendre**
4 **polynomials and B-splines**

5 Mehdi Momen¹, Malachy T. Campbell¹, Harkamal Walia², and Gota
6 Morota^{1*}

7 ¹Department of Animal and Poultry Sciences, Virginia Polytechnic Institute
8 and State University, Blacksburg, VA, USA 24061

9 ²Department of Agronomy and Horticulture, University of
10 Nebraska-Lincoln, Lincoln, NE, USA 68583

11 Keywords: genomic prediction, phenomics, longitudinal modeling, random regression, time
12 series

13

14 Running title: Longitudinal genomic prediction

15

16 ORCID: 0000-0002-2562-2741 (MM), 0000-0002-8257-3595 (MTC), 0000-0002-9712-5824 (HW),
17 and 0000-0002-3567-6911 (GM).

18

19 * Corresponding author:

20

21 Gota Morota

22 Department of Animal and Poultry Sciences

23 Virginia Polytechnic Institute and State University

24 175 West Campus Drive

25 Blacksburg, Virginia 24061 USA.

26 E-mail: morota@vt.edu

27

28 Abstract

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

49 Background

50 Plant biology has become a large-scale, data-rich field with the development of high-throughput
51 technologies for genomics and phenomics. This has increased the feasibility of data driven ap-
52 proaches to be applied to address the challenge of developing climate-resilient crops (Tester
53 and Langridge, 2010). Crop responses to environmental changes are highly dynamic and
54 have a strong temporal component. Such responses are also known as function-valued traits,
55 for which means and covariances along the trajectory change continuously. Single time
56 point measurements of phenotypes, however, only provide a snapshot, posing a series of
57 challenges for research efforts aimed at understanding the ability of the plant to mount a
58 tolerant response to an environmental constraint. Advancements in high-throughput phe-
59 notyping (HTP) technologies have enabled the automated collection of measurements at
60 high temporal resolution to produce high density image data that can capture a plethora of
61 morphological and physiological measurements (Furbank and Tester, 2011). In particular,
62 image-based phenotyping has been deemed a game changer because conventional phenotyp-
63 ing is laborious and often involves destructive methods, precluding repeated sampling over
64 time from the same individual (Ge et al., 2016). More importantly, these HTP systems offer
65 greater potential to uncover the time-specific molecular events driven by important genes
66 that have yet to be discovered in genome-wide association studies (GWAS) or to perform
67 forecasting of future phenotypes in longitudinal genomic prediction. Thus, integrating these
68 HTP data into quantitative genetics has the potential to increase the rate of genetic gain in
69 crops. However, to take full advantage of such opportunities, novel statistical methods that
70 can fully leverage time series HTP data need to be developed.

71 Recently, Campbell et al. (2018) used a random regression model (RRM) to perform ge-
72 nomic prediction for longitudinal HTP traits in controlled environments, such as greenhouses,
73 using Legendre polynomials as the choice of a basis function to model dependencies across
74 time. They also demonstrated that RRM could be used to achieve reasonable prediction
75 accuracy in a cross-validation (CV) framework to forecast future phenotypes based on infor-

76 mation from earlier growth stages. RRM also enables the calculation of (co)variances and
77 genetic values at any time between the beginning and end of the trajectory, even including
78 time points that lack phenotypic information. This study showed that RRM could effectively
79 describe the temporal dynamics of genetic signals by accounting for mean and covariance
80 structures that are subjected to change over time (Kirkpatrick et al., 1990). However, the
81 utility of RRM for plants under an abiotic stress environment is not explored. This is a crit-
82 ical unknown as the crop productivity is greatly limited by environmental challenges such
83 as drought and heat stress. In addition to the Legendre polynomials, spline functions can
84 be used to describe the relationships between image-based phenomics and genomics data
85 in longitudinal modeling (White et al., 1999). In particular, B-spline functions have been
86 used to study a variety of traits, such as growth records, in animal breeding in terms of
87 model goodness of fit using pedigree data (e.g., Meyer, 2005; Baldi et al., 2010); however, its
88 application to HTP data in plants and its predictive ability from a CV perspective remains
89 untested.

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

95 **Materials and Methods**

96 **Plant materials and greenhouse conditions**

97 Three hundred fifty-seven accessions ($n = 357$) of the rice (*O. Sativa*) diversity panel 1
98 (RDP1) were selected for this study (Zhao et al., 2011). Seeds were surface sterilized with
99 Thiram fungicide and germinated on moist paper towels in plastic boxes for three days. For
100 each accession, three uniformly germinated seedlings were selected and transplanted to pots
101 (150mm diameter x 200 mm height) filled with 2.5 kg of UC Mix. Square containers were
102 placed below each pot to allow water to collect. The plants were grown in saturated soil on
103 greenhouse benches prior to phenotyping.

104 All lines were genotyped with 44,000 single nucleotide polymorphisms (SNPs) (Zhao
105 et al., 2011). We used PLINK v1.9 software (Purcell et al., 2007) to remove SNPs with a
106 call rate ≤ 0.95 and a minor allele frequency ≤ 0.05 . Missing genotypes were imputed using
107 Beagle software version 3.3.2 (Browning and Browning, 2007). Finally, 36,901 SNPs were
108 retained for further analysis.

109 **Experimental design and drought treatment**

110 All experiments were conducted at the Plant Accelerator, Australian Plant Phenomics Fa-
111 cility, at the University of Adelaide, SA, Australia. The panel was phenotyped for a digital
112 metric representing shoot growth over 20 days of progressive drought using an image-based
113 phenomics platform. The details of the experimental design are provided in Campbell et al.
114 (2018). Briefly, each experiment consisted of 357 accessions from RDP1 and was repeated
115 three times from February to April 2016. Two smart-houses were used for each experiment.
116 In each smart-house, the accessions were distributed across 432 pots positioned across 24
117 lanes. The experiments followed a partially replicated paired design, where plants of the
118 same accession were grown adjacent to one another. In each experiment, 54 accessions were
119 randomly selected and replicated twice.

120 Seven days after transplant (DAT), plants were thinned to one seedling per pot. Two
121 layers of blue mesh were placed on top of the pots to reduce evaporation. The plants were
122 loaded on to the imaging system and were watered to 90% field capacity (FC) DAT. On
123 the 13 DAT, each pot was watered to 90% and was imaged to obtain an initial phenotype
124 before the onset of drought. One plant from each pair was randomly selected for drought
125 treatment. Water was withheld from drought plants until 25% FC, and after which water
126 was applied to maintain 25% FC. For the duration of the experiment, the control plants were
127 maintained at 100% FC.

128 **Statistical analysis of phenotypic data**

129 Visible images were processed, and digital features were extracted using the open-source
130 Python library Image Harvest (Knecht et al., 2016). The sum of plant pixels from the
131 two sides and one top view of red/green/blue (RGB) images was summed and used as a
132 measure of shoot biomass. This digital phenotype is referred to as the projected shoot area
133 (PSA) throughout this study. Several studies have reported a high correlation between PSA
134 estimates and shoot biomass (Campbell et al., 2015; Goltzarian et al., 2011; Knecht et al.,
135 2016). Prior to downstream analyses, outlier plants at each time point were detected for
136 each trait using the 1.5 interquartile range rule, and potential outliers were plotted along
137 with their treatment counterparts and inspected visually. Plants that exhibited abnormal
138 growth patterns were removed. In total, 221 plants were removed, leaving 2,586 plants for
139 downstream analyses.

140 Raw phenotypic measurements were adjusted for downstream genetic analyses prior to
141 fitting RRM. Best linear unbiased estimators (BLUE) were computed for each accession by
142 fitting experimental effect with three levels and replication within experiment for some of the
143 accessions. We postulated that observations at each time point follow the additive genetic
144 model (\mathcal{M}): $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \mathbf{Z}\mathbf{u} + \boldsymbol{\epsilon}$, where \mathbf{X} and \mathbf{Z} are $n' \times f$ and $n' \times n$ orders of incident
145 matrices linking observations (n') to systematic effects (f) and number of accessions (n),

146 respectively, \mathbf{y} is an $n' \times 1$ vector of observations at each time point, $\boldsymbol{\beta}$ is a $f \times 1$ vector of
147 systematic effects, \mathbf{u} is a $n \times 1$ vector of BLUE for accessions, and $\boldsymbol{\epsilon}$ is an $n' \times 1$ vector of
148 residuals with $Var(\boldsymbol{\epsilon}) = \mathbf{I}\sigma_{\epsilon}^2$, where \mathbf{I} is an identity matrix. This was followed by fitting a
149 RRM-based genomic prediction approach to predict phenotypes as described below.

150 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 over time were modeled through Legendre polynomials and B-splines of time at phenotyping. The general formula for the RRM was as follows:

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

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

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

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

where $\boldsymbol{\beta}$ is a vector of solutions for fixed regressions; \mathbf{u} is the additive genetic regression coefficients; \mathbf{pe} is the permanent environmental random regression coefficients; $\boldsymbol{\epsilon}$ is the residuals; and \mathbf{Z} and \mathbf{Q} are corresponding incident matrices. We assumed a multivariate-Gaussian distribution and the variance-covariance structure of

$$Var \begin{pmatrix} \mathbf{u} \\ \mathbf{pe} \\ \mathbf{e} \end{pmatrix} = \begin{pmatrix} \mathbf{G} \otimes \mathbf{C}_u & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I} \otimes \mathbf{C}_{pe} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \mathbf{R} \end{pmatrix},$$

159 where $\mathbf{G} = \mathbf{W}_{sc}\mathbf{W}'_{sc}/p$ is the genomic relationship matrix of VanRaden (2008), where \mathbf{W}_{sc}
160 represents a centered and standardized marker matrix and p is the number of markers; \mathbf{C}_u is
161 the covariance function between the random regression coefficients for the additive genetic
162 effect; \otimes is the Kronecker product; \mathbf{C}_{pe} is the covariance function between the random
163 regression coefficients for the permanent environmental effects; and $\mathbf{R} = \mathbf{I}_n\sigma_{\epsilon(t)}^2$ is a diagonal
164 matrix of heterogeneous residuals varying across times, where σ_{ϵ}^2 is the residual variance.

165 Choice of basis function

166 The choice of the basis function to model the shape of the longitudinal measurements is
167 critical. An ideal basis function has adequate potential to capture real patterns of changes
168 in variance along a continuous scale (time) for a given trait (Meyer and Kirkpatrick, 2005).
169 In this study, we used RRM with two basis functions, i.e., Legendre polynomials (Meyer,
170 1998) and B-splines (Meyer, 2005), to describe line-specific curves for the PSA trajectory
171 over the day of imaging.

172 **Legendre polynomials:** Applying parametric shape functions for covariates of time is

173 challenging because these covariates tend to generate high correlations among trajectories
174 (Mrode, 2014). For this reason, fitting Legendre polynomials of time at recording as covari-
175 ables is a common choice to model growth curves because these polynomials greatly reduce
176 the correlations between estimated random regression coefficients and make no prior assump-
177 tions regarding the shape of the longitudinal curve. This function has been used widely in
178 animal breeding for many years (e.g., Jamrozik and Schaeffer, 1997) and has recently been
179 used in plant breeding as well (Sun et al., 2017; Campbell et al., 2018; Marchal et al., 2019).
180 Suppose d is the order of fit or degree of the polynomial. Legendre polynomials evaluated at
181 the standardized time points were computed as $\Phi = \mathbf{M}\Lambda$, where \mathbf{M} is the t_{max} by $d + 1$ ma-
182 trix containing the polynomials of the standardized time covariate $\mathbf{M}_{k+1} = \left(\frac{2(t-t_{min})}{t_{max}-t_{min}}\right)^k - 1$
183 and Λ is the $d + 1 \times d + 1$ matrix of Legendre polynomial coefficients (Kirkpatrick et al.,
184 1990). Here, $t_{min} = 1$ and $t_{max} = 20$ because PSA was measured for 20 days. We chose the
185 same orders of polynomials for fixed, additive, and permanent environmental coefficients as
186 previously described Schaeffer (2016). We compared linear ($k = 0, \dots, K_1 = K_2 = K_3 = 1$)
187 and quadratic ($k = 0, \dots, K_1 = K_2 = K_3 = 2$) Legendre polynomials in this study. Thus,
188 the numbers of regression coefficients were $d + 1 = 2$ and $d + 1 = 3$ for linear and quadratic
189 Legendre polynomials, respectively.

190 **B-splines:** Spline functions consist of individual segments of polynomials joined at specific
191 points called knots. B-splines first require determination of the total number of knots K .
192 Although a large number of knots will increase complexity, too few knots will decrease accu-
193 racy. This basis function is reported to offer several advantages, including better numerical
194 properties compared with polynomials, especially when there are high genetic variances at
195 the extremes of the trajectory period, negative correlations between the most distant time
196 point measurements, and a small number of records, particularly at the last stage of the
197 trajectory (Meyer, 2005; Misztal, 2006). Here, we used equidistant knots, and the B-spline
198 function was computed from Cox-de Boor's recursion formula (De Boor, 2001). Given a
199 preconsidered knot sequence of time t , the covariables for B-splines of degree $d = 0$ were

200 defined by assuming values of unity for all points in a given interval or zero otherwise. For
201 the i th interval given by knots

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

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

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

204 This indicates that a B-spline of degree d is simply a function of B-splines of degree $d - 1$.
205 Note that the number of random regression coefficients depends on the number of knots and
206 order of polynomials for B-splines. In general, the number of regression coefficients is given
207 by $K = s + d - 1$ (Meyer, 2005). In this study, we fitted linear B-splines with $s = 3$ or
208 $s = 4$ knots to divide the time points into equally spaced intervals. The same number of
209 knots was considered for fixed trajectories, additive genetic, and permanent environmental
210 coefficients. Thus, the numbers of regression coefficients were three ($k = 1, \dots, K_1 = K_2 =$
211 $K_3 = 3 + 1 - 1 = 3$) and four ($k = 1, \dots, K_1 = K_2 = K_3 = 4 + 1 - 1 = 4$) for $s = 3$ and
212 $s = 4$ knots, respectively.

213 **Goodness of model fit**

214 The goodness of fit of RRM was assessed by computing the Akaike's information crite-
215 rion (AIC) (Akaike, 1974) and the Schwarz-Bayesian information criterion (BIC) (Schwarz
216 et al., 1978). The best model was selected based on the largest AIC and BIC values after
217 multiplying by $-1/2$. We used Wombat software to fit RMM in this study (Meyer, 2007).

218 Cross-validation scenarios

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

223 **CV1:** In the first CV scenario (CV1), the whole data set was divided into two subsets, i.e.,
224 training and testing sets, each including 179 and 178 accessions, respectively. All 20 time
225 points in the training set were fit to the RRM using Legendre polynomials and B-splines,
226 and we predicted phenotypic values of 20 time points for lines in the testing set. Random
227 assignment of individuals into the training and testing sets was repeated 10 times. The
228 equation for CV1 was set up in the following manner. The time-specific genetic value of the
229 i th individual in the training set was computed as $\hat{\mathbf{g}}_{\text{trn}, i}^t = \Phi_t \mathbf{u}_i$, where $\hat{\mathbf{g}}_{\text{trn}, i}^t$ is the estimated
230 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 ;
231 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
232 values of individuals in the testing set at time t was obtained as $\hat{\mathbf{g}}_{\text{tst}}^t = \mathbf{G}_{\text{tst}, \text{trn}} \mathbf{G}_{\text{trn}, \text{trn}}^{-1} \hat{\mathbf{g}}_{\text{trn}}^t$,
233 where $\mathbf{G}_{\text{tst}, \text{trn}}$ is the genomic relationship matrix between the testing and training set and
234 $\mathbf{G}_{\text{trn}, \text{trn}}^{-1}$ is the inverse of genomic relationship matrix between the training set. Because CV1
235 is not a forecasting task, a standard multi-trait model (MTM) was also fitted as a baseline
236 model considering longitudinal traits as different traits (Henderson and Quaas, 1976). The
237 BLUPF90 family of programs was used to fit MTM with 20 traits (Misztal et al., 2002).

238

239 **CV2:** The second CV scenario (CV2) was related to forecasting future phenotypes from
240 longitudinal traits at early time points. Individuals in the training set were used to fore-
241 cast their yet-to-be observed PSA values at later time points from information available at
242 earlier time points. The first quarter of the time points $\{t = 1, 2, 3, 4, 5\}$ was used as
243 the training set, and the remaining time points $\{t = 6, 7, \dots, 20\}$ were predicted for each
244 line in the training set. This was followed by sequentially increasing the number of time

245 points used to train the model so that in the last run, three quarters of the time points $\{t$
246 $= 1, 2, \dots, 15\}$ were used in the training set to forecast phenotypes at the last quarter of
247 time points $\{t = 16, 17, 18, 19, 20\}$. This CV scenario was designed to find a sufficient set
248 of earlier time points to obtain reasonable longitudinal prediction accuracy and is known
249 as walk forward validation. We set up the CV2 equation by first estimating the random
250 regression coefficient matrix \mathbf{u} using $\Phi_{1:t}$, which was computed from time point 1 to time
251 point t . The prediction of future time points t' ($t + 1 \leq t' \leq t_{\max}$) for an individual i was
252 carried out by $\hat{\mathbf{g}}^{t'} = \Phi_{t'} \mathbf{u}_i$, where $\Phi_{t'}$ is the t' th row vector of $t_{\max} - t$ by $K + 1$ matrix Φ ;
253 and \mathbf{u}_i is the i th column vector of the number of random regression coefficients by n matrix \mathbf{u} .

254

255 **CV3:** The third CV scenario (CV3) was designed to evaluate whether it was possible to
256 reduce the phenotyping frequency while still maintaining a high prediction accuracy for the
257 last quarter of observations. We used the last case in CV2 such that time points $\{t = 1,$
258 $2, \dots, 15\}$ were used for the training set to forecast the last quarter of observations $\{t =$
259 $16, 17, 18, 19, 20\}$. We then reduced the number of time points used in the training set as
260 follows: A, observations on odd days $\{t = 1, 3, \dots, 15\}$ were used; B, observations on even
261 days $\{t = 2, 4, \dots, 14\}$ were used; C, keep one and delete two consecutive time points. In
262 CV2 and CV3 scenarios, half of the individuals were randomly selected to fit the model, and
263 the analysis was repeated 10 times. If the loss of prediction accuracy was minimal, it was
264 possible to reduce the phenotyping cost. The equation for CV3 was set up in the same way
265 as that for CV2.

266 **Data availability**

267 Genotypic data regarding the rice accessions can be downloaded from the rice diversity panel
268 website (<http://www.ricediversity.org/>). Phenotypic data used herein are available in
269 Supplementary File S1.

270 Results

271 Assessing model fit

272 Figures 2A and 2B show the box plots of the original PSA and BLUE for the phenotypic
273 trajectories over the 20 days of imaging for control and water-limited conditions. The PSA
274 for control and water-limited plants diverged significantly after 10 days of initiation of the
275 drought treatment, and the accession level difference become apparent at later growth stages
276 under control conditions. Supplementary Figure 1 shows the linear or quadratic forms of
277 Legendre polynomials and three and four knot-based B-spline curves over 20 days of imaging.
278 For Legendre polynomials, intercept, linear, and quadratic coefficients are represented in
279 black, red, and green, respectively. For B-spline, knot 1, knot 2, and knot 3 are represented
280 in black, red, and green, respectively.

281 Table 1 summarizes the goodness of fits of RRM coupled with linear and quadratic Leg-
282 endre polynomials and B-spline functions in control and water-limited conditions. For the
283 Legendre polynomials, quadratic forms require more parameters to be estimated compared
284 with linear forms. Similar to observation for B-splines, the presence of a greater number
285 of knots suggested that there were more parameters to be estimated. Under control con-
286 ditions, the best goodness of fit was obtained by linear Legendre polynomials, followed by
287 linear B-splines with three knots, linear B-splines with four knots, and quadratic Legendre
288 polynomials according to AIC scores. According to BIC scores, linear Legendre polynomials,
289 followed by linear B-splines with three knots, quadratic Legendre polynomials, and linear
290 B-splines with four knots. Under water-limited conditions, the best goodness of fit was given
291 by linear Legendre polynomials, followed by linear B-splines with three knots, quadratic Leg-
292 endre polynomials, and linear B-splines with four knots for both AIC and BIC scores. The
293 number of parameters in the model varied from 26 to 40.

294 Cross-validation

295 The results from CV1 are shown in Figure 3. This CV was designed to evaluate the accu-
296 racy of predicting testing set individuals using all time points. Under control conditions,
297 MTM performed relatively better than RRM up to day 3. The prediction accuracy of RRM
298 increased subsequently and after the 10th day of imaging, the best prediction was given
299 by linear Legendre, followed by quadratic Legendre, linear B-spline with three knots, and
300 linear B-spline with four knots. Overall, RRM performed better than MTM, and linear Leg-
301 endre was the best prediction machine throughout the growth stages. Under water-limited
302 conditions, prediction accuracy was lower compared with those of control conditions. All
303 RRM delivered higher prediction than MTM except for the first two time points. Although
304 Legendre polynomials performed better than B-splines until day 7, the difference between
305 these approaches became negligible afterward.

306 Figures 4 and 5 show the accuracy of CV2 under control and water-limited conditions,
307 respectively. This CV was designed to test how much information from the previous time
308 points was required to achieve reasonable prediction accuracy at later growth stages. Under
309 control conditions, we found that the best prediction for the last five time points was achieved
310 when using all time point information up to the most recent (15/5 CV2 subscenario). This
311 suggested that having more information from previous time points to train the model would
312 result in higher prediction accuracy. Using the first five time points to train the model
313 resulted in the worse prediction (5/15 CV2 subscenario). Thus, it is likely that the prediction
314 accuracy in RRM declined because we attempted to estimate numerous parameters from only
315 five time points. Legendre polynomials yielded better and more stable prediction than B-
316 splines. We observed a similar trend under water-limited conditions; that is, using more
317 previous time points to train the model resulted in higher prediction accuracy. However, the
318 accuracy of prediction was unstable and decreased dramatically. There was no noticeable
319 difference between the Legendre polynomials and B-splines in terms of performance.

320 Figures 6 and 7 show the CV3 accuracy under control and water-limited conditions,

321 respectively. We designed this CV to evaluate whether it was possible to reduce phenotyping
322 frequency and phenotyping costs without sacrificing prediction accuracy. Under control
323 conditions, the prediction accuracy of CV3A, CV3B, and CV3C all decreased relative to the
324 benchmark scenario in CV2, where all of the first 15 time points were used for the training
325 set to forecast the last five time points. Although removing two consecutive time points did
326 not improve performance (CV3C), the prediction accuracy from phenotyping every other
327 day was still relatively high (CV3A and CV3B). In general, the linear B-splines performed
328 the best, and differences between scenarios were minimal. Under water-limited conditions,
329 we observed the same trend, but the prediction accuracy was more unstable and decreased
330 relative to control conditions. The quadratic Legendre polynomials and B-splines with four
331 knots did not perform well, possibly due to overfitting.

332 Discussion

333 Image-based automated HTP technologies offer great potential for characterizing multi-
334 faceted phenotypes at high temporal resolution. The use of HTP platforms plays a pivotal
335 role in accelerating breeding efforts by providing the temporal resolution needed for cap-
336 turing adaptive responses to environmental challenges, but the development of statistical
337 methodologies to analyze image-based function-valued phenotypes has not kept pace with
338 our ability to generate HTP data. Because phenomics and genomics landscapes for plants
339 are constantly advancing, parallel efforts are required to develop tools for integrating di-
340 verse genomic and phenomic datasets characterized by high temporal resolution in genetic
341 analysis. Rice is one of the most drought sensitive cereal crops, resulting in substantial
342 yield losses. With predictions for greater climatic shifts in the future and increasing com-
343 petition for fresh water resources, research that leverages the full potential of genomics and
344 phenomics is needed to elucidate the genetic and physiological basis of drought tolerance.
345 However, there is currently a lack of information regarding the modeling of temporal HTP
346 data.

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

357 Longitudinal prediction

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

369 We observed differences in prediction accuracy across models during early growth stages;
370 however, differences were incremental when predicting later growth stages in the CV1 sce-
371 nario, in which the training and testing sets were partitioned based on individuals. Overall,
372 linear Legendre polynomials performed the best and was clearly an advancement over the
373 MTM. Prediction performance in CV2, in which the training and testing sets were parti-
374 tioned according to growth stages rather individuals, showed that it was possible to predict
375 future phenotypes from information available from earlier trajectories. Here, linear and
376 quadratic Legendre polynomials produced the highest and most stable prediction accuracy
377 under control conditions, whereas linear B-splines with three knots performed better in the
378 water-limited environment. The final scenario (CV3) demonstrated that we could decrease
379 the phenotyping frequency by only phenotyping every other day to reduce the phenotyping
380 cost while minimizing the loss of prediction accuracy. In this case, linear B-spline with three
381 knots performed relatively well.

382 B-spline functions require two parameters (the position of the knots and the number
383 of knots) to be tuned. The position of knots can be chosen based on a trajectory pattern

384 such that more knots are placed for rapidly changing time points, whereas less knots are
385 placed for time points with slow changes (Misztal, 2006). Thus, the position of knots should
386 be carefully chosen if the number of phenotyped individuals varies substantially at each
387 growth stage. We chose equidistant knots in the current study because all accessions were
388 phenotyped on the same days during the trajectory. The number of knots determines the
389 number of segments fitted. When more knots are specified, the model becomes more complex.
390 Although we used $s = 3$ and $s = 4$ based on previous literature and a visual inspection of the
391 observed phenotypic trajectory, further investigations are warranted to explore the impact
392 of the number of knots on prediction accuracy. The performance of quadratic B-spline
393 functions was not evaluated in the current study because we encountered convergence issues,
394 possibly due to the small sample size. In general, we found that the advantages of B-splines
395 in inferential tasks compared with Legendre polynomials were not shown clearly in terms
396 of prediction. This is likely because PSA trajectories were relatively simple exponential or
397 monotonically increasing trajectories without obvious inflection points, indicating that the
398 potential of B-splines was not able to be fully exploited in the current study.

399 **Choice of parameters**

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

409 **Future perspective**

410 We anticipate that the current work will guide us to conduct genomic selection of econom-
411 ically important traits on the longitudinal scale for the purpose of breeding crops that are
412 adaptable to new environments or to less favorable challenging climatic conditions. More-
413 over, identifying genomic components over trajectories will provide information regarding
414 the optimum time points to maximize cost-effective selection or to design a genome-assisted
415 breeding program aiming to change the shape of the longitudinal response curve (Schaeffer,
416 2004). Using our approach, we could evaluate all changes in plant biomass accumulation
417 during the course of the experiment, in contrast to single time point analyses. Thus, we
418 expect that RRM analysis will become the norm for modeling trajectories of function-valued
419 HTP data because such approaches could be considered an extension of the widely used
420 genomic best linear unbiased prediction model for time series data. Lastly, the utility of the
421 RRM does not preclude its use in other applications. For example, the RRM offers a new
422 avenue for performing longitudinal GWAS (e.g., Howard et al., 2015; Campbell et al., 2019)
423 and genotype-by-environment interactions using the reaction norm (Arnold et al., 2019).
424 In summary, an RRM using Legendre polynomial or spline functions could be an effective
425 option for modeling trait trajectories of HTP data and may have potential applications in
426 characterizing phenotypic plasticity in plants.

427 **Acknowledgments**

428 This work was supported by the National Science Foundation under Grant Number 1736192
429 to HW and GM, and Virginia Polytechnic Institute and State University startup funds to
430 GM.

431 **Author contribution statement**

432 MTC and HW designed and conducted the experiments. MM analyzed the data. MM
433 and GM conceived the idea and wrote the manuscript. MTC and HW discussed results
434 and revised the manuscript. GM supervised and directed the study. All authors read and
435 approved the manuscript.

436 **Table**

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	CF	Log L	-0.5 AIC	-0.5 BIC	<i>p</i>
CON	LEGL	-32414.493	-32440.493	-32529.839	26
	LEGQ	-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
	LEGQ	-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; CF: covariance function; LEGL: Legendre polynomial linear; LEGQ: 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.

437 **Figures**

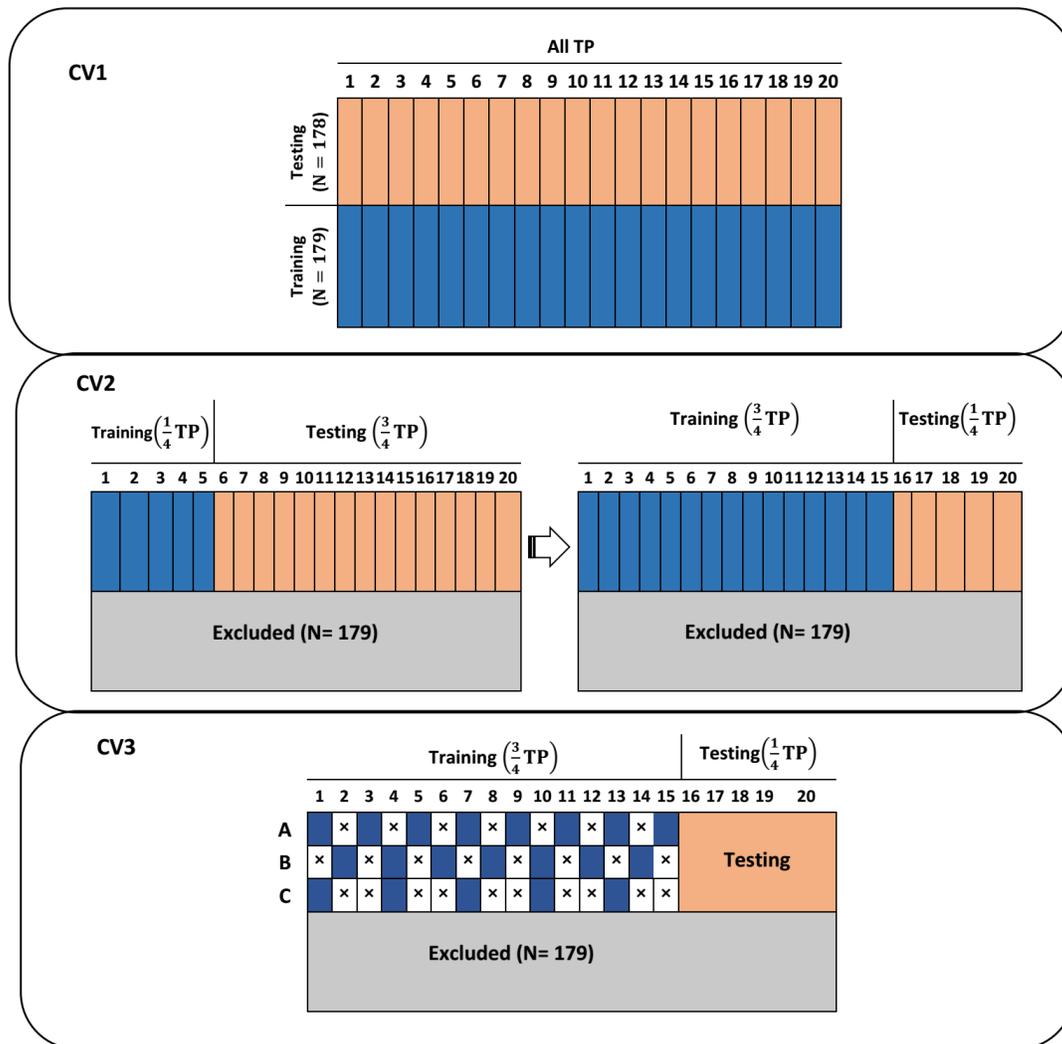


Figure 1: Pictorial representation of three cross-validation schemes used for predicting longitudinal projected shoot area (PSA) using a random regression model coupled with Legendre polynomials and B-splines. The data set consisted of 357 lines. CV1: 179 lines were used as the training set to predict 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.

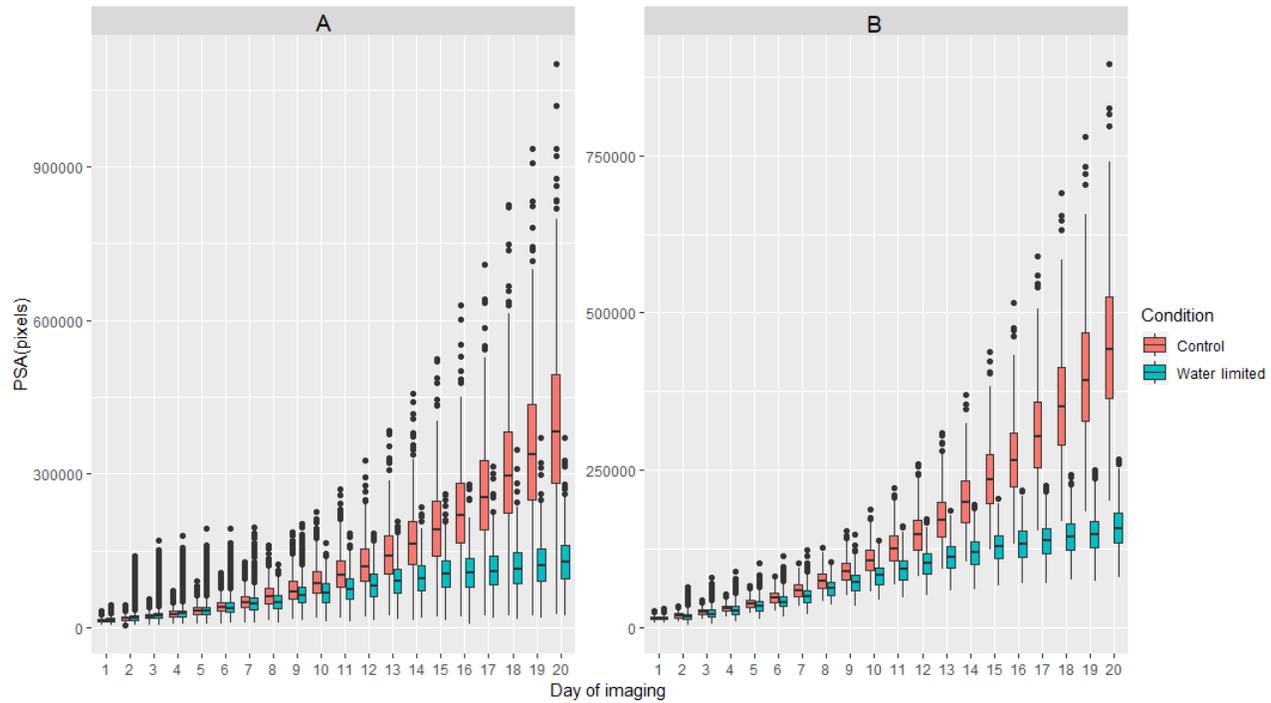


Figure 2: A: Box plots of projected shoot area (PSA) over the 20 days of imaging in two environments: controlled and water-limited conditions. B: Best linear unbiased estimators over the 20 days of imaging in two environments: controlled and water-limited conditions.

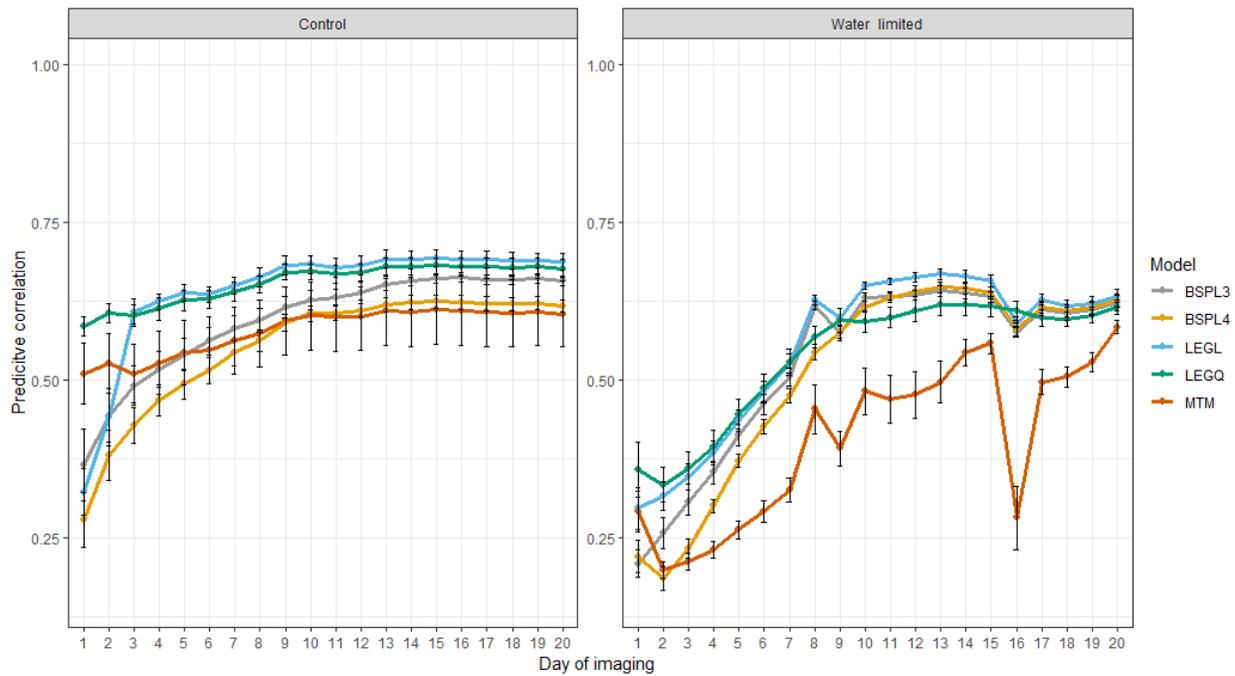
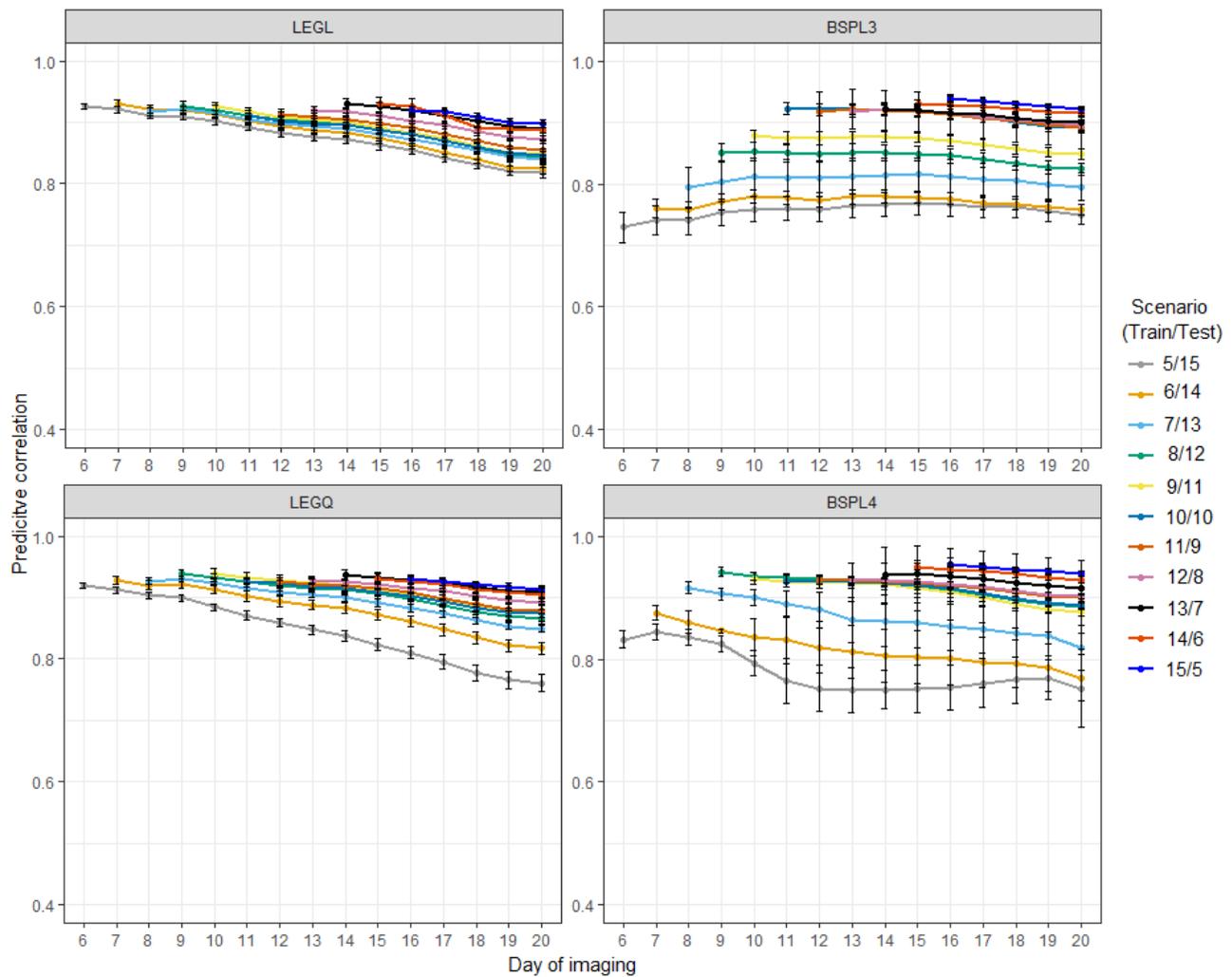


Figure 3: Prediction accuracy obtained from cross-validation 1 scenario. Total of 179 lines were used as the training set to predict PSA for the remaining 178 lines. Here, all 20 time points in the training set were used to predict PSA at each of 20 time points for a new set of lines. LEGL: linear Legendre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



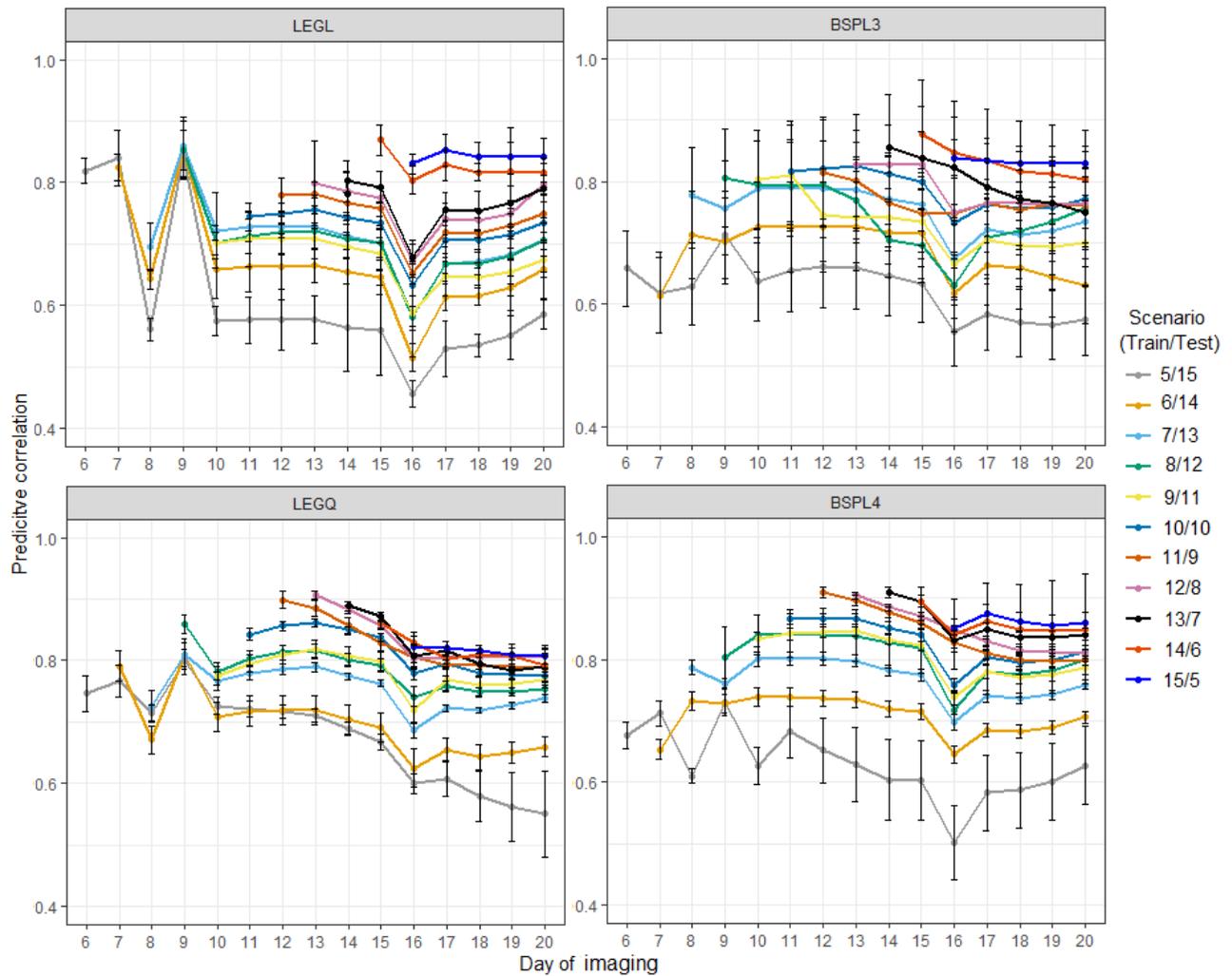


Figure 5: Prediction accuracy of cross-validation scenario 2 in water-limited conditions. Each line depicts the different number of training and testing sets partitioning at the time point levels. LEGL: linear Legendre polynomials; LEGQ: quadratic Legendre polynomials; BSPL3: linear B-splines with three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

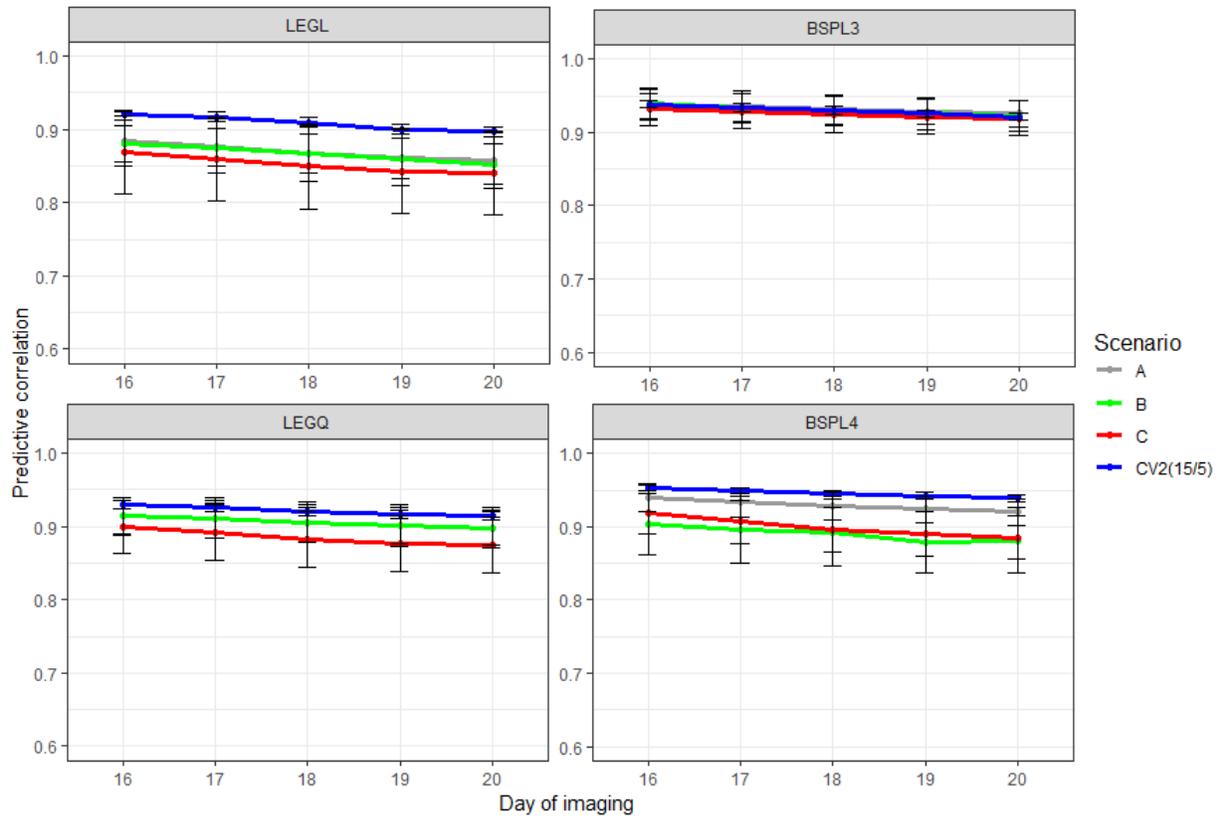


Figure 6: Prediction accuracy of cross-validation scenario 3 in control conditions. A: only observations in the odd days were used; B: only observations in the even days were used; C: keep one and delete two consecutive time points; CV2: use all available previous time points; LEGL: linear Legendre polynomials;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

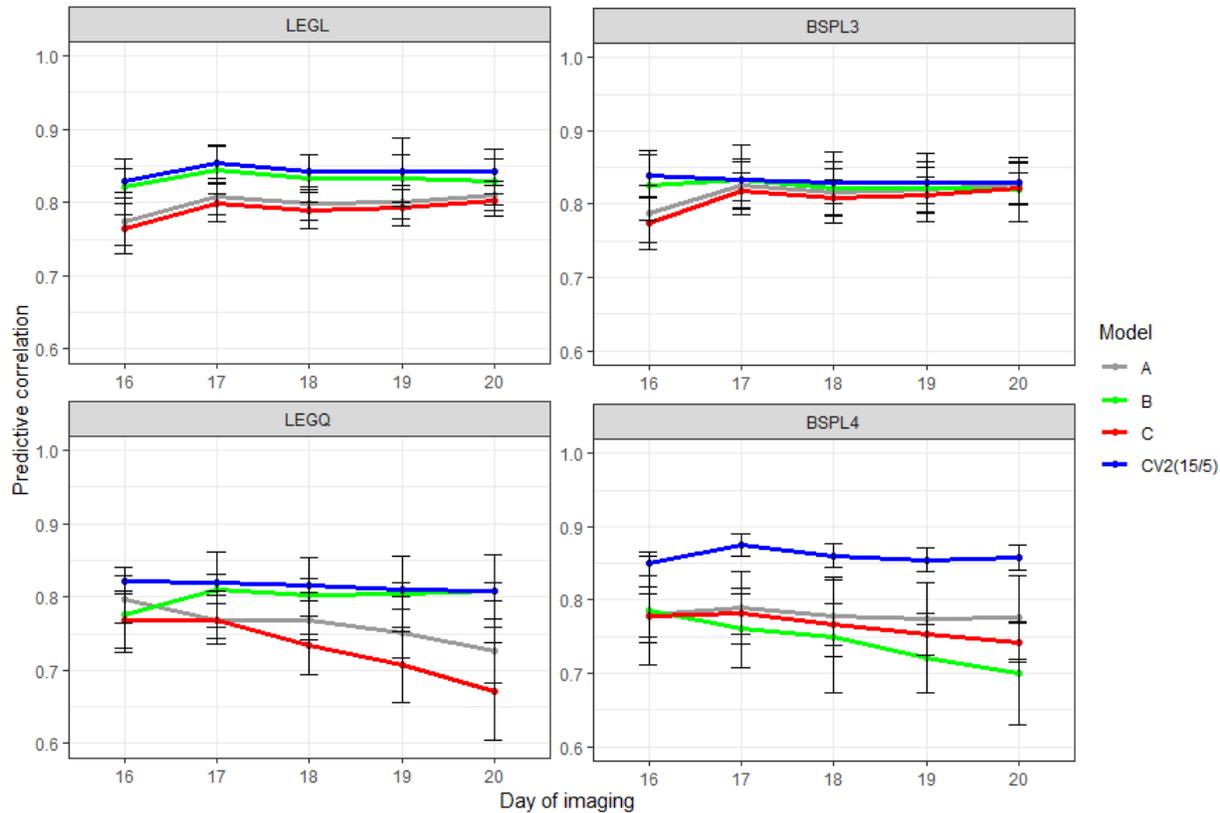


Figure 7: Prediction accuracy of cross-validation scenario 3 in water-limited conditions. A: only observations in the odd days were used; B: only observations in the even days were used; C: keep one and delete two consecutive time points; CV2: use all available previous time points; LEGL: linear Legendre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

438 References

- 439 Akaike, H. (1974). A new look at the statistical model identification. In *Selected Papers of*
440 *Hirotsugu Akaike*, pages 215–222. Springer.
- 441 Arnold, P. A., Kruuk, L. E., and Nicotra, A. B. (2019). How to analyse plant phenotypic
442 plasticity in response to a changing climate. *New Phytologist*.
- 443 Baldi, F., Alencar, M., and Albuquerque, L. G. d. (2010). Random regression analyses using
444 b-splines functions to model growth from birth to adult age in canchim cattle. *Journal of*
445 *Animal Breeding and Genetics*, 127(6):433–441.
- 446 Browning, S. R. and Browning, B. L. (2007). Rapid and accurate haplotype phasing and
447 missing-data inference for whole-genome association studies by use of localized haplotype
448 clustering. *The American Journal of Human Genetics*, 81(5):1084–1097.
- 449 Campbell, M. T., Knecht, A. C., Berger, B., Brien, C. J., Wang, D., and Walia, H. (2015).
450 Integrating image-based phenomics and association analysis to dissect the genetic archi-
451 tecture of temporal salinity responses in rice. *Plant Physiology*, 168(4):1476–1489.
- 452 Campbell, M. T., Walia, H., and Morota, G. (2018). Utilizing random regression models
453 for genomic prediction of a longitudinal trait derived from high-throughput phenotyping.
454 *Plant Direct*, 2(9).
- 455 Campbell, M. T., Walia, H., and Morota, G. (2019). Leveraging breeding values obtained
456 from random regression models for genetic inference of longitudinal traits. *The Plant*
457 *Genome*, page Early view.
- 458 De Boor, C. (2001). *A Practical Guide to Splines (Revised Edition)*, volume 27. Springer-
459 Verlag New York.
- 460 Furbank, R. T. and Tester, M. (2011). Phenomics-technologies to relieve the phenotyping
461 bottleneck. *Trends Plant Sci.*, 16:635–644.

- 462 Ge, Y., Bai, G., Stoerger, V., and Schnable, J. C. (2016). Temporal dynamics of maize
463 plant growth, water use, and leaf water content using automated high throughput rgb and
464 hyperspectral imaging. *Computers and Electronics in Agriculture*, 127:625–632.
- 465 Golzarian, M. R., Frick, R. A., Rajendran, K., Berger, B., Roy, S., Tester, M., and Lun,
466 D. S. (2011). Accurate inference of shoot biomass from high-throughput images of cereal
467 plants. *Plant Methods*, 7(1):2.
- 468 Henderson, C. and Quaas, R. (1976). Multiple trait evaluation using relatives' records.
469 *Journal of Animal Science*, 43(6):1188–1197.
- 470 Howard, J. T., Jiao, S., Tiezzi, F., Huang, Y., Gray, K. A., and Maltecca, C. (2015). Genome-
471 wide association study on legendre random regression coefficients for the growth and feed
472 intake trajectory on duroc boars. *BMC Genomics*, 16:59.
- 473 Iwaisaki, H., Tsuruta, S., Misztal, I., and Bertrand, J. (2005). Genetic parameters estimated
474 with multitrait and linear spline-random regression models using gelbvieh early growth
475 data. *Journal of Animal Science*, 83(4):757–763.
- 476 Jamrozik, J. and Schaeffer, L. (1997). Estimates of genetic parameters for a test day model
477 with random regressions for yield traits of first lactation holsteins. *Journal of Dairy
478 Science*, 80(4):762–770.
- 479 Kirkpatrick, M., Lofsvold, D., and Bulmer, M. (1990). Analysis of the inheritance, selection
480 and evolution of growth trajectories. *Genetics*, 124(4):979–993.
- 481 Knecht, A. C., Campbell, M. T., Caprez, A., Swanson, D. R., and Walia, H. (2016). Image
482 Harvest: an open-source platform for high-throughput plant image processing and analysis.
483 *Journal of Experimental Botany*, 67(11):3587–3599.
- 484 Marchal, A., Schlichting, C. D., Gobin, R., Balandier, P., Millier, F., Muñoz, F., Pâques,

- 485 L. E., and Sánchez, L. (2019). Deciphering hybrid larch reaction norms using random
486 regression. *G3: Genes, Genomes, Genetics*, 9(1):21–32.
- 487 Meyer, K. (1998). Estimating covariance functions for longitudinal data using a random
488 regression model. *Genetics Selection Evolution*, 30(3):221.
- 489 Meyer, K. (2005). Random regression analyses using B-splines to model growth of australian
490 angus cattle. *Genetics Selection Evolution*, 37(6):473.
- 491 Meyer, K. (2007). Wombat - A tool for mixed model analyses in quantitative genetics by re-
492 stricted maximum likelihood (reml). *Journal of Zhejiang University Science B*, 8(11):815–
493 821.
- 494 Meyer, K. and Hill, W. B. (1997). Estimation of genetic and phenotypic covariance functions
495 for longitudinal or repeated records by restricted maximum likelihood. *Livest Prod Sci.*,
496 47:185–200.
- 497 Meyer, K. and Kirkpatrick, M. (2005). Up hill, down dale: quantitative genetics of cur-
498 vaceous traits. *Philosophical Transactions of the Royal Society B: Biological Sciences*,
499 360(1459):1443–1455.
- 500 Misztal, I. (2006). Properties of random regression models using linear splines. *Journal of*
501 *Animal Breeding and Genetics*, 123(2):74–80.
- 502 Misztal, I., Tsuruta, S., Strabel, T., Auvray, B., Druet, T., Lee, D., et al. (2002). Blupf90
503 and related programs (bgf90). In *Proceedings of the 7th World Congress on Genetics*
504 *Applied to Livestock Production*, volume 33, pages 743–744.
- 505 Mrode, R. A. (2014). *Linear models for the prediction of animal breeding values*. CABI.
- 506 Purcell, S., Neale, B., Todd-Brown, K., Thomas, L., Ferreira, M. A., Bender, D., Maller,
507 J., Sklar, P., De Bakker, P. I., Daly, M. J., et al. (2007). Plink: a tool set for whole-

- 508 genome association and population-based linkage analyses. *The American Journal of*
509 *Human Genetics*, 81(3):559–575.
- 510 Schaeffer, L. (2016). *Random regression models*. Available in [http://animalbiosciences.](http://animalbiosciences.uoguelph.ca/~lrs/BOOKS/rrmbook.pdf)
511 [uoguelph. ca/~ lrs/BOOKS/rrmbook. pdf](http://animalbiosciences.uoguelph.ca/~lrs/BOOKS/rrmbook.pdf).
- 512 Schaeffer, L. and Jamrozik, J. (2008). Random regression models: a longitudinal perspective.
513 *Journal of Animal Breeding and Genetics*, 125(3):145–146.
- 514 Schaeffer, L. R. (2004). Application of random regression models in animal breeding. *Livest*
515 *Prod Sci.*, 86:35–45.
- 516 Schwarz, G. et al. (1978). Estimating the dimension of a model. *The Annals of Statistics*,
517 6(2):461–464.
- 518 Sun, J., Rutkoski, J. E., Poland, J. A., Crossa, J., Jannink, J.-L., and Sorrells, M. E.
519 (2017). Multitrait, random regression, or simple repeatability model in high-throughput
520 phenotyping data improve genomic prediction for wheat grain yield. *The Plant Genome*,
521 10(2).
- 522 Tester, M. and Langridge, P. (2010). Breeding technologies to increase crop production in a
523 changing world. *Science*, 327:818–822.
- 524 Van der Werf, J., Goddard, M., and Meyer, K. (1998). The use of covariance functions and
525 random regressions for genetic evaluation of milk production based on test day records.
526 *Journal of Dairy Science*, 81(12):3300–3308.
- 527 VanRaden, P. M. (2008). Efficient methods to compute genomic predictions. *Journal of*
528 *Dairy Science*, 91(11):4414–4423.
- 529 White, I., Thompson, R., and Brotherstone, S. (1999). Genetic and environmental smoothing
530 of lactation curves with cubic splines. *Journal of Dairy Science*, 82(3):632–638.

531 Zhao, K., Tung, C.-W., Eizenga, G. C., Wright, M. H., Ali, M. L., Price, A. H., Norton,
532 G. J., Islam, M. R., Reynolds, A., Mezey, J., et al. (2011). Genome-wide association
533 mapping reveals a rich genetic architecture of complex traits in *Oryza sativa*. *Nature*
534 *Communications*, 2:467.