

1 **A simple and parsimonious generalised additive model for predicting wheat yield in a decision**
2 **support tool**

3 Kefei Chen^{1,2}, Rebecca A. O’Leary¹, Fiona H. Evans^{1,2,3,†}

4

5 ¹ Department of Primary Industries and Regional Development, 3 Baron-Hay Court, South Perth, WA
6 6151, Australia

7 ² Faculty of Science and Engineering, Curtin University, Bentley, WA 6102, Australia

8 ³ Big Data in Agriculture, Murdoch University, Murdoch, WA 6150, Australia

9

10 † Corresponding author

11 Phone: +61 (0) 427 871 776

12 Fax: +61 8 9368 3082

13 Email: Fiona.Evans@murdoch.edu.au

14

16 **Abstract**

17 Yield prediction is a major determinant of many management decisions for crop production. Farmers
18 and their advisors want user-friendly decision support tools for predicting yield. Simulation models
19 can be used to accurately predict yield, but they are complex and difficult to parameterise. The goal of
20 this study is to build a simple and parsimonious model for predicting wheat yields that can be
21 implemented in a decision tool to be used by farmers at a paddock level.

22

23 A large yield data set accumulated from trials on commonly grown varieties in Western Australia is
24 used to build and validate a generalised additive model (GAM) for predicting wheat yield.

25 Explanatory variables tested included weather data and derivatives, geolocation, soil type, land
26 capability, and wheat varieties. Model selection followed a forward stepwise approach in combination
27 with cross-validation to select the smallest set of explanatory variables. The predictive performance is
28 also evaluated using independent data.

29

30 The final model uses seasonal water availability, location and year to predict wheat yield. Because the
31 GAM model has minimal inputs, it can be easily employed in a decision tool to predict yield
32 throughout the growing season using rainfall data up to the prediction date and either climatological
33 averages or seasonal forecasts of rainfall for the remainder of the growing season. It also has the
34 potential to be used as an input to agronomic models that predict the effect on yield of various
35 management choices for fertilizer, pest, weed and disease management.

36

37 **Keywords:** Yield prediction; Waterlogging; Precision farming; Crop modelling; Crop water relations;
38 Decision support

39

40

41 **1. Introduction**

42 Grain production in Western Australia (WA) is located in the south-west of Australia and wheat is the
43 main crop grown. The climate is Mediterranean, with broadacre cropping reliant on winter rainfall in
44 a dryland system. Climate variability explains around 40% of total wheat yield variability in the
45 Australian wheat belt, and in parts of WA that figure can be greater than 60% (Ray et al., 2015).

46

47 The WA grainbelt has experienced a 20% decline in winter rainfall since the 1970s due to southward
48 shifts in rain-bearing synoptic systems (Bates et al., 2008; Hope et al., 2006). Despite the challenge of
49 declining rainfall, grain growers have improved productivity; largely through better use of existing
50 technologies, including decision support tools and precision agriculture technologies (Kingwell et al.,
51 2013).

52

53 Decision support tools help farmers assess different tactical management choices by predicting their
54 effects on crop yield and profit. Many incorporate agronomic models that apply Mitscherlich's law of
55 diminishing returns, which states that the increase in crop yield as a result of increasing a single
56 growth factor is proportional to the decrease from the maximum yield that takes into account all other
57 limiting factors (Mitscherlich, 1909). For instance, this approach is used in decision support tools for
58 nitrogen application in common use in WA such as 'NPdecide', 'Select Your Nitrogen' and its
59 derivative the N-Broadacre mobile application (Bowden, 2003; Burgess et al., 1991). The user of
60 these decision tools is required to input an estimate of maximum yield when nitrogen is not a limiting
61 factor. In the absence of better information, many farmers and / or agricultural advisors using decision
62 tools use an average estimate of yield. Given that wheat yields can vary considerably from year to
63 year depending on the amount, frequency and timing of seasonal rainfall, the benefits of using
64 decision support tools for on-farm management can be greatly improved with better prediction of
65 yield to use in the tools.

66

67 With further reduction in rainfall likely for WA (IOCI, 2012), improvements in the performance of
68 decision support tools via better yield prediction can help grain growers maintain and / or continue to
69 improve wheat yields despite a drying climate.

70

71 In addition, yield predictions can be applied regionally to support strategic decision-making and
72 planning processes in agribusiness in the context of the increasing pressure on food security and
73 sustainability. Predictions of crop yield before harvest can assist transport and storage planning,
74 allowing grain handling organisations to best allocate resources in any given season. Crop yield
75 predictions have been used as an aid in land use planning, determining crop insurance premiums
76 (Abbaspour et al., 1992; Choudhury and Jones, 2014)-and food supply considerations (Rosengrant et
77 al., 2002), and considering future production in a changing climate (Ludwig and Asseng, 2006).

78 Regional yield predictions can also be used to understand where and why gaps between potential and
79 actual yield occur (van Ittersum et al., 2013).

80

81 Prediction of crop yield is usually by one of two methods: (1) crop simulation modelling or (2)
82 statistical crop modelling. Crop simulation models, such as APSIM (Ahmed et al., 2016; Brown et al.,
83 2018; Keating et al., 2003; McCown et al., 1996; Yang et al., 2014), CERES (Lv et al., 2017; Ritchie
84 et al., 1988), SWAP (Mokhtari et al., 2018; van Lier et al., 2015) and WOFOST (Boogaard et al.,
85 1998; Ma et al., 2013), use fundamental mechanisms of crop growth and development, soil and water
86 processes to simulate plant growth in different scenarios. They are mathematical representations of
87 the real-world situation and usually include sub-models for crop growth, soil water movement,
88 fertiliser uptake and dissolution and more. Crop simulation models give accurate predictions of yield,
89 but require extensive parameter inputs that describing characteristics of the modelled situation
90 (including crop, soil, climate and crop management). Use of crop simulation models for on-farm
91 decision making typical requires soil testing and meticulous calibration, often with the assistance of a
92 specialised consultant (Hunt et al., 2006). For regional yield predictions, a generic set of typical model
93 parameters are used, but these parameters will give inaccurate predictions for any individual location
94 and are therefore not of use in decision support tools.

95

96 Statistical crop models encode relationships between environmental and/or management factors and
97 crop yield to make predictions. The most commonly used methods in WA are French and Schultz
98 (F&S) type equations that relate the total amount of water available during the growing season with
99 the potential non-water-limited yield possible in Australian dryland cropping systems (French and
100 Schultz, 1984). They developed the following equation: $Potential\ yield\ (kg/ha) = (Wavail -$
101 $Evaporation) \times WUE$, where *Wavail* is the amount of water that is available during the growing
102 season, defined as $Wavail = RF_S/3 + RF_{GS}$, where RF_S is summer rainfall (November-March) and
103 RF_{GS} is growing season rainfall (April-October). French and Schultz estimated values of 110 mm for
104 *Evaporation* and 20 kg/ha/mm for the water use efficiency (*WUE*).

105

106 A recent modification frequently used in WA is the broken stick method that specifies an upper limit
107 to potential yield dependent on the plant available water capacity of the soil (Oliver et al., 2009). This
108 has been shown to better account for waterlogging, which has adverse effects on the development and
109 growth of wheat in Mediterranean-type environments (Turner, 1992). Both the F&S and broken stick
110 methods have the advantage of being simple, easy to understand, and require minimal inputs and have
111 been employed in decision support tools such as PYCAL (Tennant and Tennant, 2000), CliMate
112 (<https://climateapp.net.au>) and the Department of Primary Industries and Regional Development
113 (DPIRD) 'Potential yield tool' (<https://www.agric.wa.gov.au/climate-weather/potential-yield-tool>).
114 However, these equations can be perceived as being too simple and they only provide a prediction of
115 potential non-water-limited yield not actual yield.

116

117 Outside of WA, the most common statistical method used for predicting crop yields is multiple linear
118 regression (linear models) which has been used with covariates including seasonal climate data
119 (Gornott and Wechsung, 2016; Landau et al., 2000; Lobell and Burke, 2010), remote sensing
120 vegetation indices (Kern et al., 2018; Qader et al., 2018; Zhang et al., 2017) and crop water stress
121 indices (Schut et al., 2009), using either simulated data from crop models or real data to train the
122 models. Bayesian implementations of linear models have also been applied (Bornn and Zidek, 2012;

123 Chipanshi et al., 2015), and partial least squares regression has been used to identify patterns in daily
124 climate records that have use for predicting yield (Ceglar et al., 2016). Machine learning approaches
125 to yield prediction from seasonal climate data include artificial neural network (ANN) models (Çakır
126 et al., 2014; Das et al., 2018); random forests, which have been used at global and / or regional scales
127 (Folberth et al., 2019; Jeong et al., 2016); and genetic programming algorithms (Ali et al., 2018).

128

129 The use of time series models, including auto-regressive integrated moving average (ARIMA)
130 models, has been investigated for forecasting crop yields (Ali et al., 2015; Choudhury and Jones,
131 2014; Craparo et al., 2015; Debnath et al., 2013). However, these models rely on the assumption that
132 yield in any one year is related to that of the previous year and because wheat yields in WA vary
133 considerably from year to year with seasonal rainfall, they are unsuitable for this study.

134

135 Motivated by the need for seasonal yield predictions for use in decision support tools for farmers, the
136 objective of this study is to determine whether a large yield data set can be used to build a statistical
137 model that combines the advantages of simple F&S type equations with added complexity, as
138 required, to produce results akin to those of crop simulation models. The goal is to produce a model
139 that is more accurate than simple linear equations while still having few, easily obtained inputs so that
140 it can be implemented in a decision tool. For this reason, we limit climate inputs to seasonal
141 summaries rather than daily data, and compare model outputs only with F&S and not with crop
142 simulation models that require daily climate data as input.

143

144 A generalized additive model (GAM) is a generalized linear model in which the predictor depends
145 linearly on unknown smooth functions of some predictor variables, and interest focuses on inference
146 about these smooth functions. GAMs were originally developed to blend properties of generalized
147 linear models with additive models (Hastie and Tibshirani, 1990). GAMs allow for flexible
148 specification of the dependence of the response (yield) on covariates, by specifying the model in
149 terms of smooth functions or ‘smooths’, rather than as detailed parametric relationships (Wood,
150 2006). GAMs have a number of advantages that suggest their use for crop yield modelling – most

151 importantly that they capture nonlinear relationships as compared to multivariate linear regression.
152 GAMs are easy to interpret by plotting fitted smooths against yield. Relationships between smooths
153 and yield can be considered in light of existing knowledge to ensure that they are sensible. They can
154 also identify hidden patterns in the data, potentially improving our understanding of crop-soil-weather
155 interactions.

156

157 This study uses a large yield data set accumulated from variety trials on commonly grown and well-
158 adapted wheat varieties in WA. The data encompass a wide range of locations, soils and rainfall
159 zones. Explanatory variables tested in the study were selected using existing knowledge and following
160 consultation with local agronomists. They include weather data and derivatives, geolocation (latitude
161 and longitude), soil type, differences in wheat variety and various classifications of wheat variety (eg.
162 into classes of maturity type).

163

164 Forward stepwise model selection is used to build a GAM that has the smallest possible set of inputs
165 required to predict yield. The forward selection process starts with the null model (i.e. intercept only
166 model) and adds variables to the model one at a time. At each step, the ability of each additional
167 explanatory variable to add predictive performance to the model is assessed using statistics for model
168 selection that balance the ‘goodness of fit’ of the model with its complexity. The variable that
169 provides the most improvement is then added to the model, until no further improvement is possible.
170 At each step, cross-validation is applied to test the ability of the model to perform on unseen data and
171 thus avoid overfitting. The final model is also tested using previously unseen data from a separate
172 source.

173

174 This approach is used to build a simple and parsimonious GAM that reflects current understanding of
175 how a wheat crop responds to seasonal rainfall, and how yields vary spatially and through time in
176 WA. Because the GAM has minimal inputs, it can be easily employed in a decision tool to predict
177 yield throughout the growing season using rainfall data up to the prediction date and either
178 climatological averages or seasonal forecasts of rainfall for the remainder of the growing season. It

179 also has the potential to be used as an input to agronomic models that predict the effect on yield of
180 various management choices for farm inputs, pest, weed and disease management.

181

182 **2. Data**

183 *2.1. Yield data*

184 *Variety trial data:* A total of 25,505 observations from variety trials, including WA-based Crop
185 Variety Trials (CVT) and National Variety Trials (NVT, <http://www.nvtonline.com.au>) in a 40-year
186 period of 1975 – 2014 were considered in this study. The data include 109 varieties and 775 unique
187 locations; however, the majority of varieties are not commonly grown and/or not well-adapted in
188 Western Australia (WA). A panel of 26 varieties commonly grown and well-adapted in WA with a
189 total of 17,701 observations were selected for the base model construction. Figure 1 shows the
190 distribution of 775 locations in variety trials conducted in the WA grainbelt, overlaid on 30-year
191 average rainfall (Garlinge 2005).

192

193 *Focus paddock data:* A total of 428 observations from 164 different Focus Paddocks (Harries et al.,
194 2015) over the period of 2010 – 2013 were used to test the predictive performance of the base model
195 constructed using the variety trial data. Unlike the variety trial data which was specifically collected
196 to determine performance differences between wheat varieties, the Focus Paddocks included data
197 from operational wheat-growing farm paddocks. They were mainly distributed in the medium rainfall
198 zones (see Figure S1) and were used in this study as an independent test of the fitted yield prediction
199 model to help understand how the model might predict previously unseen, real paddock data.

200

201 *2.2. Weather and derived data*

202 Weather data for the nearest weather station to each trial location were extracted from the Patched
203 Point Database (<https://www.longpaddock.qld.gov.au/silo>). Data derived from the daily weather
204 include:

205

206 *Growing season available water (Wavail)*: The amount of water available for crop use in any growing
207 season is defined to be one third of summer rainfall plus growing season rainfall.

208 *30-year average rainfall*: Average rainfall for years 1971 to 2000 (*avgrf30*) was used as a continuous
209 surrogate for rainfall zone.

210

211 *Early season rainfall*: Early season rainfall (*sumrf2m*) is defined to be the sum of the first two months
212 of rainfall occurring after germination.

213

214 *Germination time*: Germination time (*gdoy*) is estimated from daily rainfall using a common
215 germination rule where germination is assumed to occur if there is 25mm of rainfall over three days
216 after 25 April, or 5mm of rainfall over three days after 5 June.

217

218 *2.3. Soil data*

219 *Soil type*: The dominant soil types for each location were extracted from the soils database managed
220 by the Department of Primary Industries and Regional Development (DPIRD). The soil groups were
221 first classified according to their soil-landscape mapping (Purdie et al., 2004), and then simplified to
222 six general functional types (*soils*) including clays, duplex, gravel, loamy, sandy and wet to reduce the
223 number of categories for modelling. The classification of the agricultural soil groups to functionally
224 simplified soil classes is shown in Table S1.

225

226 *Land capability and qualities*: A total of 19 variables related to land capability and qualities were
227 derived from the Natural resource information (NRInfo) database maintained by DPIRD and other
228 government agencies (<https://www.agric.wa.gov.au/resource-assessment/nrinfo-western-australia>).

229 The *in-silico* land capability and qualities include variables of land evaluation for percentage of map
230 unit for soil pH value, phosphorus export hazard, surface salinity, drainage potential, subsurface
231 compaction susceptibility, water erosion hazard, waterlogging susceptibility, water repellence
232 susceptibility, soil water storage, wind erosion hazard, and dryland cropping (van Gool et al., 2005).

233 A detailed description of the 19 variables related to the land capability and qualities is presented in
234 Table S2.

235

236 Other explanatory variables of interest for the yield prediction model were also assessed, including
237 wheat variety (*variety*) and maturity type (*mattype*).

238

239 **3. Methodology**

240 *3.1. Modelling using GAMs*

241 A generalised additive modelling (GAM) approach (Hastie and Tibshirani, 1990) is used for the yield
242 prediction. In particular, we used the 'gam' function from R package 'mgcv' (Wood, 2006).

243

244 *3.2. Model selection*

245 A forward stepwise model selection is used to determine the optimal set of explanatory variables to
246 use in the model. The forward selection process starts with the null model (i.e. intercept only model)
247 and adds variables to the model one at a time. At each step, variables that provide the most
248 information to the model are calculated using various statistics for model selection.

249

250 *3.3. Statistics for model performance*

251 Several model performance statistics were used for model selection and to assess the goodness-of-fit
252 of the models including the Akaike information criterion (AIC) (Akaike, 1973), Bayesian information
253 criterion (BIC) (Schwarz, 1978), root mean square error (RMSE), correlation coefficient (r) and
254 coefficient of determination (R^2).

255

256 *3.4. Cross validation*

257 When fitting statistical models, the goodness-of-fit of the model must be balanced against model
258 complexity in order to avoid overfitting. That is, to avoid building models that describe the data used
259 to fit them very well but predict poorly on previously unseen data. To this end, the 'gam' function in

260 R package `'mgcv'` can be used to fit smooth terms by generalised cross validation (GCV). GCV is
261 based on a leave-one-out cross validation, where only one datum from the dataset is omitted from
262 model fitting (Craven and Wahba, 1978). For small data sets, this results in reasonable smooths;
263 however, for a large data set such as ours, this resulted in smooths that appeared to be over-fitted. We
264 therefore applied a 5-fold-cross-validation (Geisser, 1993), to divide the variety trial data into 5
265 testing and training sets with a 80:20 split of training to test GAMs at each step of the forwards model
266 selection procedure.

267

268 *3.5. Model checking with deviance residuals*

269 The final fitted GAM was checked with deviance residuals. A variety of residual plots were
270 examined, including normal Q-Q plot, residuals versus fitted values, histogram of residuals and
271 response versus fitted values.

272

273 *3.6. Bayesian analysis*

274 Bayesian inference using Gibbs sampling approach was used for comparison to the GAM approach.
275 The `'jagam'` function in the `'mgcv'` package in R was used to convert the best GAM model from
276 the frequentist approach into a Bayesian `'JAGS'` (Just Another Gibbs Sampler) model, following the
277 approach outlined in Wood (2016). Applying the `'jagam'` function converts the GAM into a
278 Bayesian graphical model for simulation with `'JAGS'`, which can then be passed to `'JAGS'` via the
279 `'rjags'` package in R. The following model diagnostics were applied to test the convergence of the
280 Markov chain Monte Carlo (MCMC) simulations: Geweke's Z-score (Geweke, 1991), Gelman and
281 Rubin's convergence diagnostic (Gelman and Rubin, 1992) and Heidelberger and Welch's
282 convergence diagnostic (Heidelberger and Welch, 1983). In addition, the convergence was also
283 checked using a variety of plots, including a trace density plot, Gelman-Rubin-Brooks plot and
284 Geweke-Brooks plot.

285

286 **4. Results**

287 4.1. Exploratory data investigation

288 In the variety trial data, there are a large number of observations with the same location, variety and
289 year, but with large within-cluster variation (see Figure S2). To reduce the within-cluster variance
290 component, the data were aggregated over location, variety, and year. Thus, the average yield within
291 the same location, year and variety was used for model construction. After aggregation, there were a
292 total of 9,116 observations for the variety trial data set.

293

294 Figure 2 shows changes in mean yield and *Wavail* through time. The mean yield and *Wavail* with
295 standard error of the mean (SEM) are shown in Table S3. There were obvious fluctuations of the
296 mean yield and *Wavail* through the 40-year period, and the overall trends of the mean yield and
297 *Wavail* differ markedly. The correlation between yield and *Wavail* in the early years of the variety
298 trial data is much lower compared to that in the late years. This suggests that seasonal rainfall has
299 played a more important role in crop production as the *Wavail* has gradually declined in WA.

300 Regardless of the declining trend of *Wavail*, yield has gradually increased over these years. This could
301 possibly be attributed to the improvement in crop management techniques and technologies, including
302 crop breeding for better selection of varieties and timely sowing. For this reason, we have tested the
303 year in which the trial was conducted as a potential variable in the model.

304

305 Figure 3 shows the distribution of yield, and its relationship with explanatory variables *Wavail*,
306 *sumrf2m*, *avgrf30*, germination time (*gdoy*), longitude (*longi*), latitude (*lat*), year, variety and soils.
307 Because the data are from variety trials, a range of different varieties are planted simultaneously at
308 each trial location causing a wide range in recorded yields shown in the y-axes. The correlation of
309 yield with each continuous variable is marked on the plots in Figure 3, showing that seasonal
310 available water, *Wavail*, has the highest correlation with yield.

311

312 As may be expected, rainfall in the first two months after germination and 30-year average rainfall are
313 correlated with *Wavail*, with correlations of 0.73 and 0.69 respectively. Both show similar, but
314 weaker, relationships with yield as for *Wavail*. Germination time (*gdoy*) has a significant correlation

315 with yield ($p\text{-value} < 2.2 \times 10^{-16}$) (Figure 3), which reflects regional knowledge that longer
316 growing periods (i.e. earlier sowing) results in higher yields in the absence of other constraints such as
317 frost damage.

318

319 Analysis of variance shows that inter-variety variation in yield is significant ($p\text{-value} < 2.2 \times 10^{-16}$);
320 however, the amount of variability in yield explained by variety is only 2.5%. The same is true for
321 soils. Similarly, the five maturity types explained less than 0.5% of variability in yield.

322

323 4.2. Model selection

324 A selected list of models for water-limited yield with the best performance at each level is presented
325 in Table 1, with the model formulae as specified in R. A full list of models for yield compared in this
326 study, including model performance statistics, is presented as supplementary material Table S4.

327 Seasonal water availability (W_{avail}) was selected as the first variable added by stepwise forward
328 selection using a linear model. It had the greatest coefficient of determination (R^2) compared to all of
329 the univariate analyses, showing that more variation can be explained by W_{avail} than any of the other
330 explanatory variables. This supports our prior knowledge that seasonal rainfall is the main driver of
331 yield in water-limited environments. Use of a GAM with the single explanatory variable W_{avail} ,
332 showed that using a smooth term provided a better fit and resulted in a lower AIC and RMSE than the
333 linear model, and higher correlation coefficient and R^2 . Thus, the GAM with W_{avail} as a smooth (s)
334 term was selected as a base model. The degrees of freedom of the smooth are dictated by the
335 dimension, k , of the spline basis used. We optimised k between numbers of 3 to 10 using 5-fold cross-
336 validation. According to the test statistics and examination of the smooth plots for W_{avail} , increasing
337 k beyond 3 was not required because higher values only increased the complexity of the smooth while
338 resulting in greater values of the AIC (see Table S4). Therefore, after the first step of forward model
339 selection, the selected model had the form: $yield \sim s(W_{avail}, k=3)$ (see Figure 4). The plot of the fitted
340 smooth against yield shows that the partial effect of W_{avail} on yield reaches a maximum around 400
341 mm, after which it starts to have a negative effect on yield (Figure 5a).

342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368

The second step in the forward model selection tested the addition of smooth terms for each explanatory variable, as well as a two-dimensional smooth term fitted to capture the spatial correlated effect of geolocation, $s(lat, longi)$. Adding the smooth term $s(lat, longi)$ to the base model provided the greatest improvement, and it was therefore selected as the second covariate to add into the model. Five-fold cross validation showed that setting k equal to 6 was optimal. The contour plot for $s(lat, longi)$ is similar to the grainbelt rainfall contours, and is a realistic representation of prior knowledge that yield decreases with both rainfall and distance from the Indian and Southern oceans (Figure 5b). After adding the second term, the model after performing the step 2 of forward selection is $yield \sim s(Wavail, k=3) + s(lat, longi, k=6)$.

Following the same process, the third covariate to be added to the model was the variable year (Table 1). There was no benefit gained by adding year as a smooth term, and therefore year was fitted in model as a linear effect.

Adding soil types to the model could provide slightly better model performance based on 5-fold cross-validation (Figure 4), but the model does not inherently possess higher predictive power based on a test with out-of-sample data (i.e. Focus Paddocks data). Since this model is built for a predictive purpose, a simple and parsimonious model with few inputs was selected as the final fitted model: $yield \sim s(Wavail, k=3) + s(lat, longi, k=6) + year$.

4.3. Model checking with deviance residuals

The residual plots (Figure S3) indicate that the GAM model adequately fits the data, but there is a large amount of variability in the observed data that is not captured by the model. Because this model is not accounting for any kind of farm or paddock management (e.g. fertiliser applications, soil applications or weed, disease and pest management), this variability is expected.

369 *4.4. Spatial and temporal model diagnostics*

370 Figure 6 shows the GAM model predicted yields compared to the observed yield for each year in the
371 variety trial data. The model clearly performs better in later years, most obviously 2000. This is likely
372 due to the change in the relationship between *Wavail* and yield (see Figure 2). Figure S4 confirms
373 this, and also shows that model performance is better in low to medium seasonal rainfall.

374

375 *4.5. Model performance using unseen data*

376 The Focus Paddocks data were used as an independent test of the fitted yield prediction model. They
377 were not used to fit the model. The correlation between observed and predicted yields was 0.549 and
378 the RMSE was 0.948 tonnes/ha. Compared with the results of the 5-fold cross-validation, which
379 showed a correlation of 0.655 and RMSE of 0.833 tonnes/ha calculate over the test sets. Figure 7
380 shows the annual plots of observed versus predicted yield for the Focus Paddock data. The GAM
381 model tends to under-predict yield for the Focus Paddocks, particularly in years 2011 and 2012. These
382 years are both examples of types of years that have already been identified as years in which the
383 model performs poorly: in 2011 there was a poor correlation between yield and *Wavail* in the variety
384 trial data; and 2012 was a wetter year.

385

386 The plot for the smoothing term $s(Wavail)$ in the GAM model is displayed over the Focus Paddock
387 data in Figure S5a. This suggests the effect of *Wavail* on yield is approximately linear over this data
388 set. Similarly the smoothing term $s(lat, longi)$ using Focus Paddocks data is displayed in Figure S5b.

389

390 *4.6. Bayesian analysis*

391 The selected “best” GAM model was used to create a Bayesian(‘JAGS’) model code and data, which
392 was then run using ‘JAGS’ to make inference about the model via Gibbs sampling. The trace and
393 density plot, Gelman-Rubin-Brooks plot, and the Geweke-Brooks plot for all parameters in the Gibbs
394 simulations are presented in Figure S6, Figure S7 and Figure S8, respectively. The Geweke
395 convergence diagnostic suggests that there is a good equity of the means of the first (10%) and last
396 part (50%) of a Markov chain. The Gelman-Rubin convergence diagnostic indicates that the output

397 from all chains is indistinguishable. The Heidelberger-Welch convergence diagnostic suggests that
398 there is little evidence to reject the null hypothesis i.e. the Gibbs sampled values come from a
399 stationary distribution. The results of these convergence diagnostics are presented in supplementary
400 Table S5. Together, these diagnostics indicate that the fitted Bayesian 'JAGS' model adequately
401 describes the variety trial data.

402

403 The resulting 'JAGS' predictions were similar to those of the original GAM model with ($r = 0.999$),
404 suggesting that no benefit is gained by use of a full 'JAGS' implementation of this model.

405

406 *4.7. Comparison with French & Schultz method*

407 The comparisons between the (traditional) GAM predicted yield and French & Schultz (F&S)
408 potential yield using variety trial data and Focus Paddock data are presented in Figure S9 and Figure
409 S10, respectively. The potential yield was estimated using the F&S approach with 120 mm as
410 evaporation and 20 kg/ha.mm as the water use efficiency. In general, predicted yields are lower than
411 potential yields, as would be expected. When *Wavail* is low, the GAM predicted yield is slightly
412 higher than the potential yield estimated by F&S method (Figure S9d).

413

414 **5. Discussion**

415 The model selection process used in this study ensured that goodness-of-fit was balanced against
416 model complexity by adding explanatory variables one at a time, only if they added to the predictive
417 ability of the model. This process showed that water availability is the major determinant of wheat
418 yield in Western Australia. Based on a historical trial series from 1975 to 2014, the partial effect of
419 *Wavail* on yield reaches a plateau around 400 mm, and then it starts to have a negative effect on yield
420 (Figure 5a).

421

422 The second most important predictor is location. The partial effect of the combination of latitude and
423 longitude is similar to the contour map of the rainfall zones. The addition of location to the model

424 adds to its predictive ability, indicating that the spatial effect is caused by more than rainfall patterns.

425 The visualised contour plot of the effect of latitude and longitude on the GAM yield prediction is
426 shown in Figure 8, in which the partial effect of latitudes and longitudes on yield are classified into
427 six different contours i.e. 1.5, 2, 2.5, 3, 3.5 and 4 (see Figure 8).

428

429 The third most important predictor is the year in which the crop is sown. This predictor is most likely
430 capturing the increasing ability of WA grain growers to grow higher yielding crops despite decreasing
431 rainfall, either by new technology or improved knowledge. *Adding soil and variety covariates did not
432 improve yield predictions, so the final fitted GAM predicts a mean yield across all varieties and soil
433 types.*

434

435 The simple and parsimonious GAM model developed in this study can be used as a base model for
436 predicting yield that can potentially be integrated with agronomic models for predicting further effects
437 of farm management for nutrition, pest, disease and weed management. The GAM model can be used
438 to predict yields throughout the growing season using rainfall data up to the prediction date and either
439 climatological averages or seasonal forecasts of rainfall for the remainder of the growing season,
440 similar to the ‘Rainfall to Date’ tool that predicts growing season rainfall , developed by DPIRD in
441 WA . The tool provides cumulative growing seasonal rainfall and projected seasonal rainfall finishes
442 at a given date based on climatological history ([https://www.agric.wa.gov.au/climate-weather/rainfall-
443 date](https://www.agric.wa.gov.au/climate-weather/rainfall-date)).

444

445 Model assessment showed better performance in more recent years, most obviously from 2000, that is
446 likely due to changes in the relationship between available water and wheat yield. The model also
447 performs better in years with low to medium rainfall, with a tendency to under-predict in high rainfall
448 years. This pattern was supported by model validation using the independent Focus Paddocks data set.

449

450 Model validation suggests that the effect of *Wavail* on yield is approximately linear (Figure S5a)
451 across the domain of the Focus Paddocks data. This is probably because the Focus Paddocks data are

452 mainly distributed in the medium rainfall zone, with few of the high *Wavail* values evident in the
453 variety trial data.

454

455 Comparison of the GAM predictions with the French & Schultz (F&S) potential yield showed that
456 GAM predicted yields were generally lower than the F&S potential yield (Figure S9d and Figure
457 S10d), as is expected because F&S aims to predict non-water-limited yield potential whereas the
458 GAM is aiming to predict actual yield taking into account other constraints.

459

460 **6. Conclusions**

461 Variation in wheat yield in Western Australia can be described using a simple GAM model with
462 minimal inputs derived from rainfall, latitude/longitude and year. This model can easily be developed
463 into a user-friendly online tool for use by grain growers and their consultants. The model can be used
464 to predict yields throughout the growing season using rainfall data up to the prediction date and either
465 climatological averages (Hunt et al., 2006), or seasonal forecasts of rainfall as used by (Brown et al.,
466 2018) for the remainder of the growing season. It also has the potential to be used as an input to
467 agronomic models that predict the effect on yield of various management choices for farm inputs,
468 pest, weed and disease management.

469

471 **7. Acknowledgments**

472 We are grateful to many colleagues from the Department of Primary Industries and Regional
473 Development (DPIRD) and grain growers for providing various data and help. This project is made
474 possible by Royalties for Regions. The authors would like to thank Dennis Van Gool, Karen Holmes,
475 Ted Griffin, Tony Leeming, Avril Russell-Brown and Phil Goulding from the GIS soil group at
476 DPIRD for providing the soil type and land capability data; many other colleagues at DPIRD for
477 helpful discussions and suggestions.

478

479 **8. Conflicts of interest**

480 The authors declare no conflicts of interest.

482 **9. References**

- 483 Abbaspour, K., Hall, J., Moon, D., 1992. A yield model for use in determining crop insurance
484 premiums. *Agricultural and Forest Meteorology* 60, 33-51.
- 485 Ahmed, M., Akram, M.N., Asim, M., Aslam, M., Hassan, F.-u., Higgins, S., Stöckle, C.O.,
486 Hoogenboom, G., 2016. Calibration and validation of APSIM-Wheat and CERES-Wheat for spring
487 wheat under rainfed conditions: Models evaluation and application. *Computers and Electronics in*
488 *Agriculture* 123, 384-401.
- 489 Akaike, H., 1973. Information theory and an extension of the maximum likelihood principle. Second
490 International Symposium on Information Theory, 267-281.
- 491 Ali, M., Deo, R.C., Downs, N.J., Maraseni, T., 2018. Cotton yield prediction with Markov Chain
492 Monte Carlo-based simulation model integrated with genetic programming algorithm: A new hybrid
493 copula-driven approach. *Agricultural and Forest Meteorology* 263, 428-448.
- 494 Ali, S., Badar, N., Fatima, H., 2015. Forecasting production and yield of sugarcane and cotton crops
495 of Pakistan for 2013-2030. *Sarhad Journal of Agriculture* 31, 1-10.
- 496 Bates, B.C., Hope, P., Ryan, B., Smith, I., Charles, S., 2008. Key findings from the Indian Ocean
497 Climate Initiative and their impact on policy development in Australia. *Climatic Change* 89, 339-354.
- 498 Boogaard, H., Van Diepen, C., Rotter, R., Cabrera, J., Van Laar, H., 1998. WOFOST 7.1; user's guide
499 for the WOFOST 7.1 crop growth simulation model and WOFOST Control Center 1.5. SC-DLO.
- 500 Bornn, L., Zidek, J.V., 2012. Efficient stabilization of crop yield prediction in the Canadian Prairies.
501 *Agricultural and forest meteorology* 152, 223-232.
- 502 Bowden, J.W., 2003. Select your Nitrogen. A decision tool for quantifying nitrogen availability and
503 crop response in broad-acre farming systems. . WA Department of Agriculture Bulletin.
- 504 Brown, J.N., Hochman, Z., Holzworth, D., Horan, H., 2018. Seasonal climate forecasts provide more
505 definitive and accurate crop yield predictions. *Agricultural and Forest Meteorology* 260, 247-254.
- 506 Burgess, S.J., Bowden, J.W., Diggle, A.J., 1991. NPDECIDE - User's Guide. A computer program to
507 help with nitrogen and phosphorus decisions for cereals. Department of Agriculture, Government of
508 Western Australia.
- 509 Çakır, Y., Kırıcı, M., Güneş, E.O., 2014. Yield prediction of wheat in south-east region of Turkey by
510 using artificial neural networks, *Agro-geoinformatics (Agro-geoinformatics 2014)*, Third International
511 Conference on. IEEE, pp. 1-4.
- 512 Ceglar, A., Toreti, A., Lecerf, R., Van der Velde, M., Dentener, F., 2016. Impact of meteorological
513 drivers on regional inter-annual crop yield variability in France. *Agricultural and forest meteorology*
514 216, 58-67.
- 515 Chipanshi, A., Zhang, Y., Kouadio, L., Newlands, N., Davidson, A., Hill, H., Warren, R., Qian, B.,
516 Daneshfar, B., Bedard, F., 2015. Evaluation of the Integrated Canadian Crop Yield Forecaster
517 (ICCYF) model for in-season prediction of crop yield across the Canadian agricultural landscape.
518 *Agricultural and Forest Meteorology* 206, 137-150.
- 519 Choudhury, A., Jones, J., 2014. Crop yield prediction using time series models. *Journal of Economics*
520 *and Economic Education Research* 15, 53.

- 521 Craparo, A., Van Asten, P., Läderach, P., Jassogne, L., Grab, S., 2015. *Coffea arabica* yields decline
522 in Tanzania due to climate change: Global implications. *Agricultural and Forest Meteorology* 207, 1-
523 10.
- 524 Craven, P., Wahba, G., 1978. Smoothing noisy data with spline functions. *Numerische mathematik*
525 31, 377-403.
- 526 Das, B., Nair, B., Reddy, V.K., Venkatesh, P., 2018. Evaluation of multiple linear, neural network and
527 penalised regression models for prediction of rice yield based on weather parameters for west coast of
528 India. *Int J Biometeorol* 62, 1809-1822.
- 529 Debnath, M., Bera, K., Mishra, P., 2013. Forecasting area, production and yield of cotton in India
530 using ARIMA model. *Research & Reviews: Journal of Space Science & Technology* 2, 16-20.
- 531 Folberth, C., Baklanov, A., Balkovič, J., Skalský, R., Khabarov, N., Obersteiner, M., 2019. Spatio-
532 temporal downscaling of gridded crop model yield estimates based on machine learning. *Agricultural*
533 *and Forest Meteorology* 264, 1-15.
- 534 French, R.J., Schultz, J.E., 1984. Water use efficiency of wheat in a Mediterranean-type environment.
535 I. The relation between yield, water use and climate. *Australian Journal of Agricultural Research* 35,
536 743-764.
- 537 Geisser, S., 1993. *Predictive inference*. Chapman and Hall, New York, NY.
- 538 Gelman, A., Rubin, D.B., 1992. Inference from iterative simulation using multiple sequences.
539 *Statistical science*, 457-472.
- 540 Geweke, J., 1991. Evaluating the accuracy of sampling-based approaches to the calculation of
541 posterior moments. Federal Reserve Bank of Minneapolis, Research Department Minneapolis, MN,
542 USA.
- 543 Gornott, C., Wechsung, F., 2016. Statistical regression models for assessing climate impacts on crop
544 yields: A validation study for winter wheat and silage maize in Germany. *Agricultural and Forest*
545 *Meteorology* 217, 89-100.
- 546 Harries, M., Anderson, G.C., Hüberli, D., 2015. Crop sequences in Western Australia: what are they
547 and are they sustainable? Findings of a four-year survey. *Crop and Pasture Science* 66, 634-647.
- 548 Hastie, T.J., Tibshirani, R.J., 1990. *Generalized additive models*. CRC Press, Boca Raton, London,
549 New York, Washington, D.C.
- 550 Heidelberg, P., Welch, P.D., 1983. Simulation run length control in the presence of an initial
551 transient. *Operations research* 31, 1109-1144.
- 552 Hope, P.K., Drosowsky, W., Nicholls, N., 2006. Shifts in the synoptic systems influencing southwest
553 Western Australia. *Climate Dynamics* 26, 751-764.
- 554 Hunt, J., van Rees, H., Hochman, Z., Carberry, P.S., Holzworth, D., Dalgliesh, N.P., Brennan, L.E.,
555 Poutlon, P.L., van Rees, S., Huth, N.I., Peake, A., 2006. Yield Prophet®: An online crop simulation
556 service, 13th Australian Agronomy Conference.
- 557 IOCI, 2012. Project 1.2 South-west Western Australia's Regional Surface Climate and Weather
558 Systems, IOCI Milestone Report

- 559 Jeong, J.H., Resop, J.P., Mueller, N.D., Fleisher, D.H., Yun, K., Butler, E.E., Timlin, D.J., Shim,
560 K.M., Gerber, J.S., Reddy, V.R., Kim, S.H., 2016. Random forests for global and regional crop yield
561 predictions. *PLoS One* 11, e0156571.
- 562 Keating, B.A., Carberry, P.S., Hammer, G.L., Probert, M.E., Robertson, M.J., Holzworth, D., Huth,
563 N.I., Hargreaves, J.N.G., Meinke, H., Hochman, Z., McLean, G., Verburg, K., Snow, V., Dimes, J.P.,
564 Silburn, M., Wang, E., Brown, S., Bristow, K.L., Asseng, S., Chapman, S., McCown, R.L., Freebairn,
565 D.M., Smith, C.J., 2003. An overview of APSIM, a model designed for farming systems
566 simulation. *European Journal of Agronomy* 18, 267-288.
- 567 Kern, A., Barcza, Z., Marjanović, H., Árendás, T., Fodor, N., Bónis, P., Bognár, P., Lichtenberger, J.,
568 2018. Statistical modelling of crop yield in Central Europe using climate data and remote sensing
569 vegetation indices. *Agricultural and forest meteorology* 260, 300-320.
- 570 Kingwell, R., Anderton, L., Islam, N., Xayavong, V., Wardell-Johnson, A., Feldman, D., Speijers, J.,
571 2013. Broadacre farmers adapting to a changing climate, in: *Facility, N.C.C.A.R. (Ed.), p. 171.*
- 572 Landau, S., Mitchell, R., Barnett, V., Colls, J., Craigon, J., Payne, R., 2000. A parsimonious, multiple-
573 regression model of wheat yield response to environment. *Agricultural and forest meteorology* 101,
574 151-166.
- 575 Lobell, D.B., Burke, M.B., 2010. On the use of statistical models to predict crop yield responses to
576 climate change. *Agricultural and Forest Meteorology* 150, 1443-1452.
- 577 Ludwig, F., Asseng, S., 2006. Climate change impacts on wheat production in a Mediterranean
578 environment in Western Australia. *Agricultural Systems* 90, 159-179.
- 579 Lv, Z., Liu, X., Cao, W., Zhu, Y., 2017. A model-based estimate of regional wheat yield gaps and
580 water use efficiency in main winter wheat production regions of china. *Scientific reports* 7, 6081.
- 581 Ma, H., Huang, J., Zhu, D., Liu, J., Su, W., Zhang, C., Fan, J., 2013. Estimating regional winter wheat
582 yield by assimilation of time series of HJ-1 CCD NDVI into WOFOST-ACRM model with Ensemble
583 Kalman Filter. *Mathematical and Computer Modelling* 58, 759-770.
- 584 McCown, R., Hammer, G., Hargreaves, J., Holzworth, D., Freebairn, D., 1996. APSIM: a novel
585 software system for model development, model testing and simulation in agricultural systems
586 research. *Agricultural systems* 50, 255-271.
- 587 Mitscherlich, E.A., 1909. Das gesetz des minimums und das Gesetz des abnehmenden Bodenertrages
588 (Eng: The law of the minimum and the law of diminishing soil productivity). *Landwirtschaftliche*
589 *Jahrbücher* 38, 537-552.
- 590 Mokhtari, A., Noory, H., Vazifedoust, M., 2018. Improving crop yield estimation by assimilating LAI
591 and inputting satellite-based surface incoming solar radiation into SWAP model. *Agricultural and*
592 *Forest Meteorology* 250, 159-170.
- 593 Oliver, Y.M., Robertson, M.J., Stone, P.J., Whitbread, A., 2009. Improving estimates of water-limited
594 yield of wheat by accounting for soil type and within-season rainfall. *Crop and Pasture Science* 60,
595 1137-1146.
- 596 Purdie, B.R., Tille, P.J., Schoknecht, N.R., 2004. Soil-landscape mapping in south-western Australia:
597 an overview of methodology and outputs., *Resource Management Technical Reports. Department of*
598 *Agriculture and Food Western Australia.*

- 599 Qader, S.H., Dash, J., Atkinson, P.M., 2018. Forecasting wheat and barley crop production in arid and
600 semi-arid regions using remotely sensed primary productivity and crop phenology: A case study in
601 Iraq. *Science of the Total Environment* 613, 250-262.
- 602 Ray, D.K., Gerber, J.S., MacDonald, G.K., West, P.C., 2015. Climate variation explains a third of
603 global crop yield variability. *Nat Commun* 6, 5989.
- 604 Ritchie, J., Godwin, D., Otter-Nacke, S., 1988. CERES-Wheat. A simulation model of wheat growth
605 and development. Univ. of Tex. Press, Austin.
- 606 Rosengrant, M., Cai, X., Cline, S.A., 2002. World water and food to 2025. International Food Policy
607 Research Institute, Washington, DC.
- 608 Schut, A.G.T., Stephens, D.J., Stovold, R.G.H., Adams, M., Craig, R.L., 2009. Improved wheat yield
609 and production forecasting with a moisture stress index, AVHRR and MODIS data. *Crop and Pasture
610 Science* 60, 60-70.
- 611 Schwarz, G., 1978. Estimating the dimension of a model. *The annals of statistics* 6, 461-464.
- 612 Tennant, D., Tennant, S., 2000. Potential Yield Calculator, software package. Department of
613 Agriculture and Food Western Australia.
- 614 Turner, N.C., 1992. Crop production on duplex soils: an introduction. *Animal Production Science* 32,
615 797-800.
- 616 van Gool, D., Tille, P.J., Moore, G.A., 2005. Land evaluation standards for land resource mapping :
617 assessing land qualities and determining land capability in south-western Australia, in: Department of
618 Agriculture and Food, W.A. (Ed.), Perth, Western Australia.
- 619 van Ittersum, M.K., Cassman, K.G., Grassini, P., Tittonell, P., Hochman, Z., 2013. Yield gap analysis
620 with local to global relevance—A review. *Field Crops Research* 143, 4-17.
- 621 van Lier, Q.d.J., Wendroth, O., van Dam, J.C., 2015. Prediction of winter wheat yield with the SWAP
622 model using pedotransfer functions: An evaluation of sensitivity, parameterization and prediction
623 accuracy. *Agricultural Water Management* 154, 29-42.
- 624 Wood, S., 2006. Generalized additive models: an introduction with R. CRC press, Boca Raton,
625 London, New York.
- 626 Wood, S.N., 2016. Just Another Gibbs Additive Modeller: Interfacing JAGS and mgcv. arXiv preprint
627 arXiv:1602.02539.
- 628 Yang, Y., Yang, Y., Han, S., Macadam, I., Li Liu, D., 2014. Prediction of cotton yield and water
629 demand under climate change and future adaptation measures. *Agricultural water management* 144,
630 42-53.
- 631 Zhang, N., Zhao, C., Quiring, S.M., Li, J., 2017. Winter Wheat Yield Prediction Using Normalized
632 Difference Vegetative Index and Agro-Climatic Parameters in Oklahoma. *Agronomy Journal* 109,
633 2700-2713.

635 **10. Tables**

636

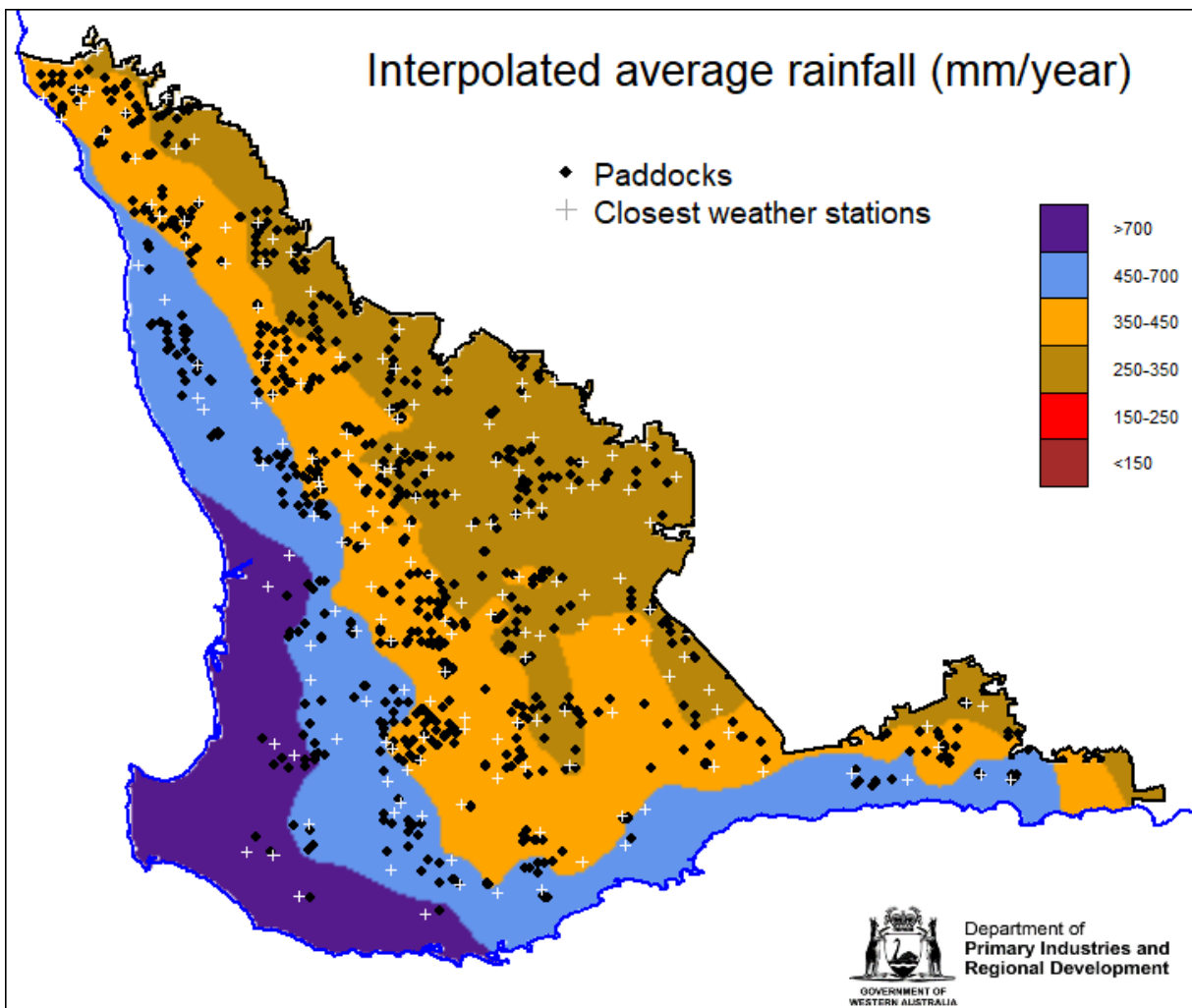
637 Table 1: Selected list of models for water-limited yield tested with cross validation in this study.

Model	Model_Formula	Level	Corr (train)	Corr (test)	RMSE (train)	RMSE (test)	R ² (train)	R ² (test)	AIC (all)	BIC (all)
1	lm(yield~Wavail)	1	0.526	0.525	0.939	0.939	0.276	0.276	24723	24744
2	gam(yield~s(Wavail, k = 3))	1	0.583	0.582	0.896	0.897	0.34	0.338	23884	23913
12	gam(yield~s(Wavail, k = 3) + s(longi, lat, k = 6))	2	0.631	0.63	0.856	0.857	0.397	0.394	23049	23113
22	gam(yield~s(Wavail, k = 3) + s(longi, lat, k = 6) + year)	3	0.656	0.655	0.832	0.833	0.43	0.427	22546	22617
24	gam(yield~s(Wavail, k = 3) + s(longi, lat, k = 6) + s(year, k = 4))	3	0.658	0.656	0.831	0.833	0.431	0.427	22522	22607
32	gam(yield~s(Wavail, k = 3) + s(longi, lat, k = 6) + year + soils)	4	0.663	0.661	0.826	0.828	0.437	0.433	22425	22532
33	gam(yield~s(Wavail, k = 3, by = soils) + s(longi, lat, k = 6) + year)	4	0.661	0.659	0.828	0.83	0.434	0.428	22473	22614

638 k refers to the degree of freedom of the smooth function; Corr: correlation coefficient;

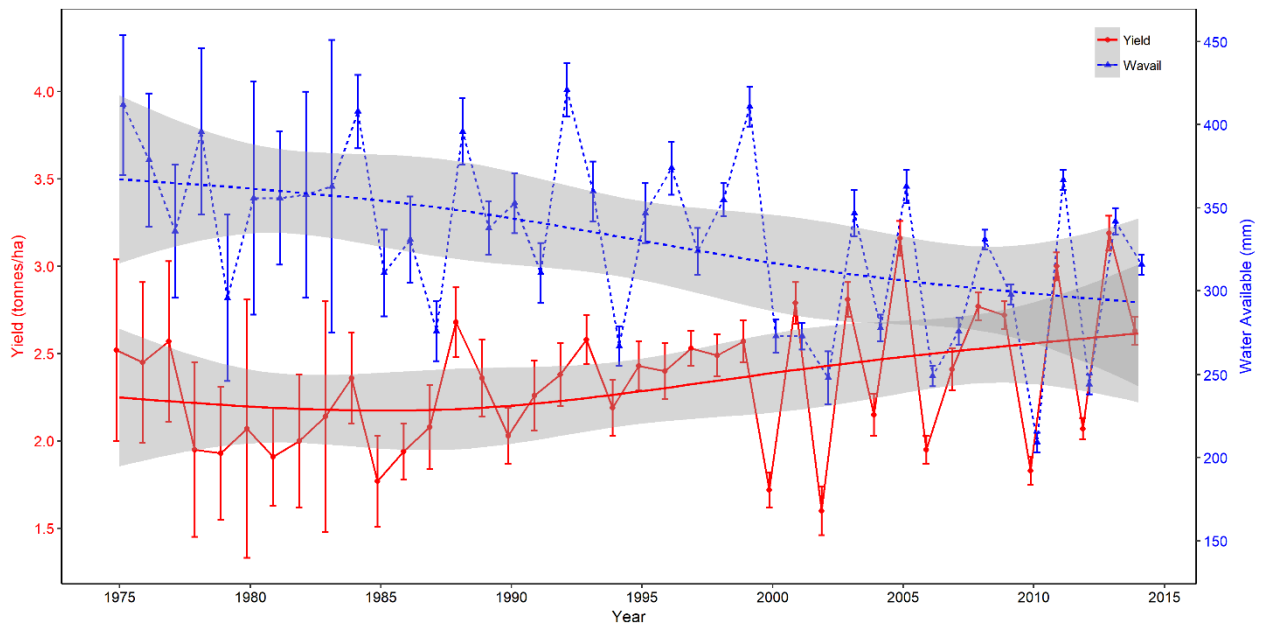
639 †: A full list of models for yield tested with cross validation in this study is provided as supplementary information (Table S4).

640



642
643 Figure 1: Map of the paddocks and patched point weather stations in variety trials with interpolation of
644 *avgrf30*.

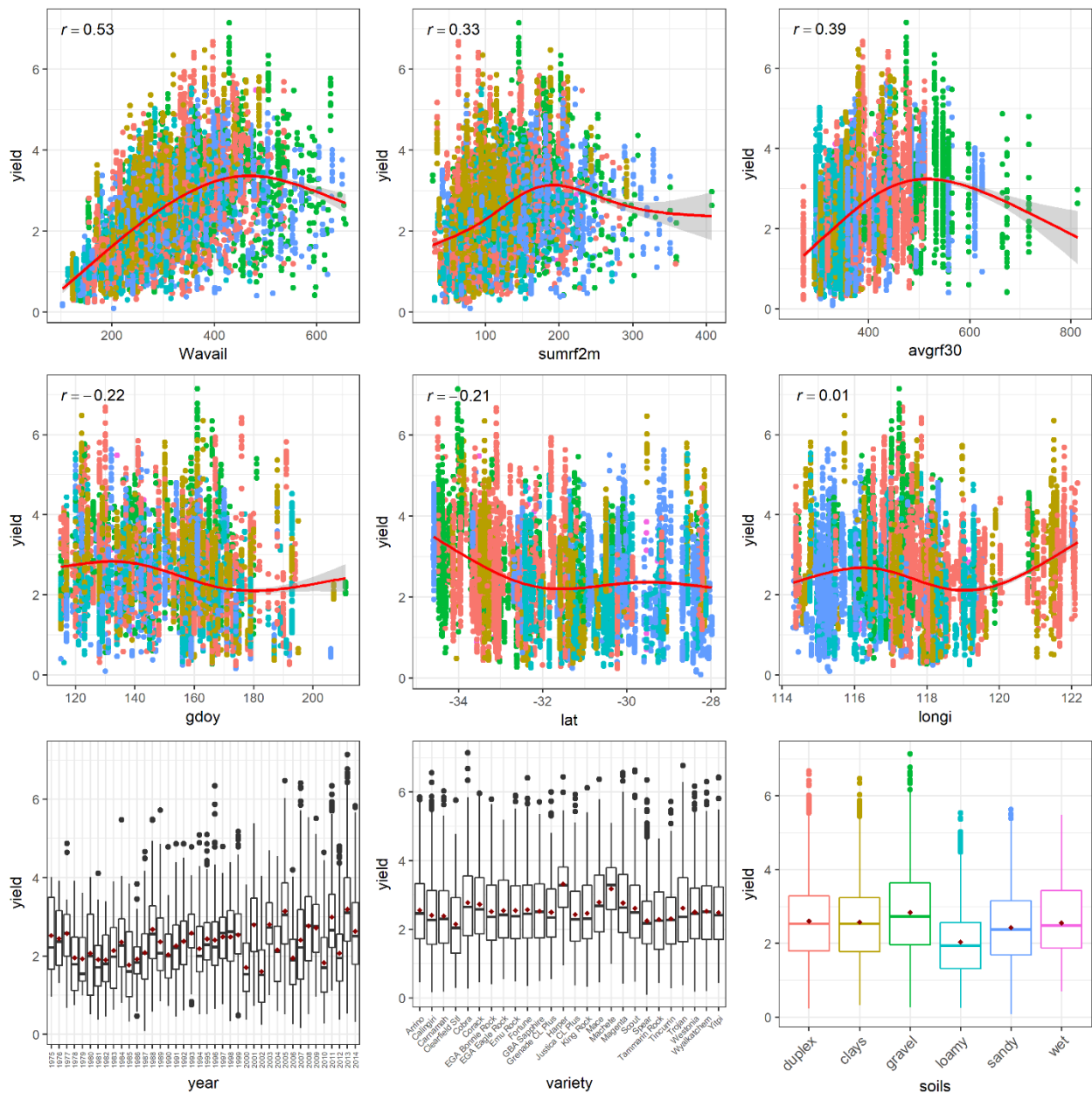
646



647

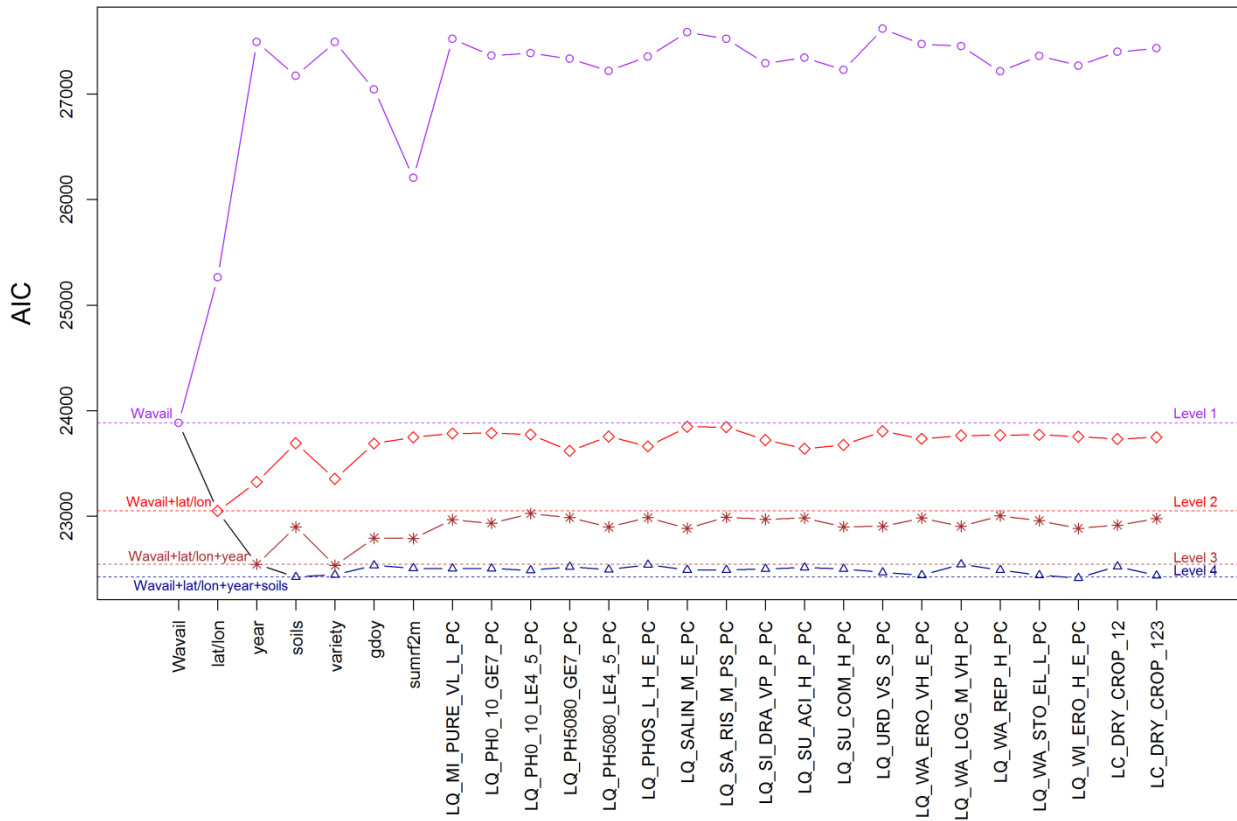
648 Figure 2: Plots of yield and *Wavail* through time. Note: The error bars represent $2 \times \text{SEM}$ (standard error of
649 the mean). The non-linear curves were fitted with natural cubic spline function with knot number of 3. The
650 shaded region represents 95% confidence intervals for the smooth curves.

651



652

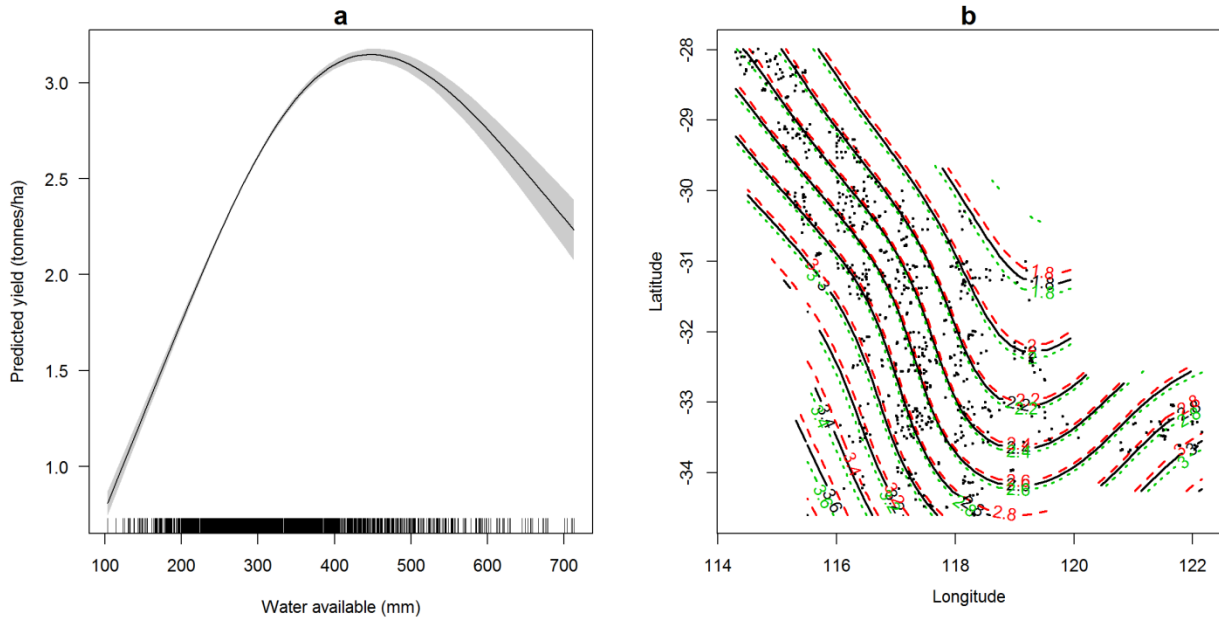
653 Figure 3: Scatter plot / boxplot matrix for the response variable (yield) and some explanatory variables
 654 including *Wavail*, *sumrf2m*, *avgrf30*, day of germination (*gdoy*), latitude (*lat*), longitude (*longi*), year,
 655 variety and soils. Scatter plots are shown for the continuous variables and boxplots for the categorical
 656 variables. The points on the scatter plots are observed yield (from the aggregated data based on year,
 657 location (latitude/longitude), variety and soil type (n=9068)) and are coloured according to soil types
 658 (red=duplex, brown=clays, green=gravel, light blue=loamy, dark blue=sandy, purple=wet). Because the
 659 data are from variety trials, a range of different varieties are planted simultaneously at each location
 660 causing a wide range in recorded yields shown in the y-axes. The non-linear curves were fitted with
 661 natural cubic spline function with knot number of 3. The shaded region on the continuous variable plots
 662 represents the 95% confidence intervals for the smooth curves. For the continuous variables, the
 663 correlation of yield with each variable is displayed in the top left-hand corner.



665

666 Figure 4: Forward stepwise model selection using AIC. Each variable that was not already in the model
 667 was tested for inclusion in the model. The most significant variable of these variables was added to the
 668 model at each level. The model started with the null-model, and followed by adding variables to the model
 669 one at a time, and continued adding variables until none of remaining variables are significant when added
 670 to the model.

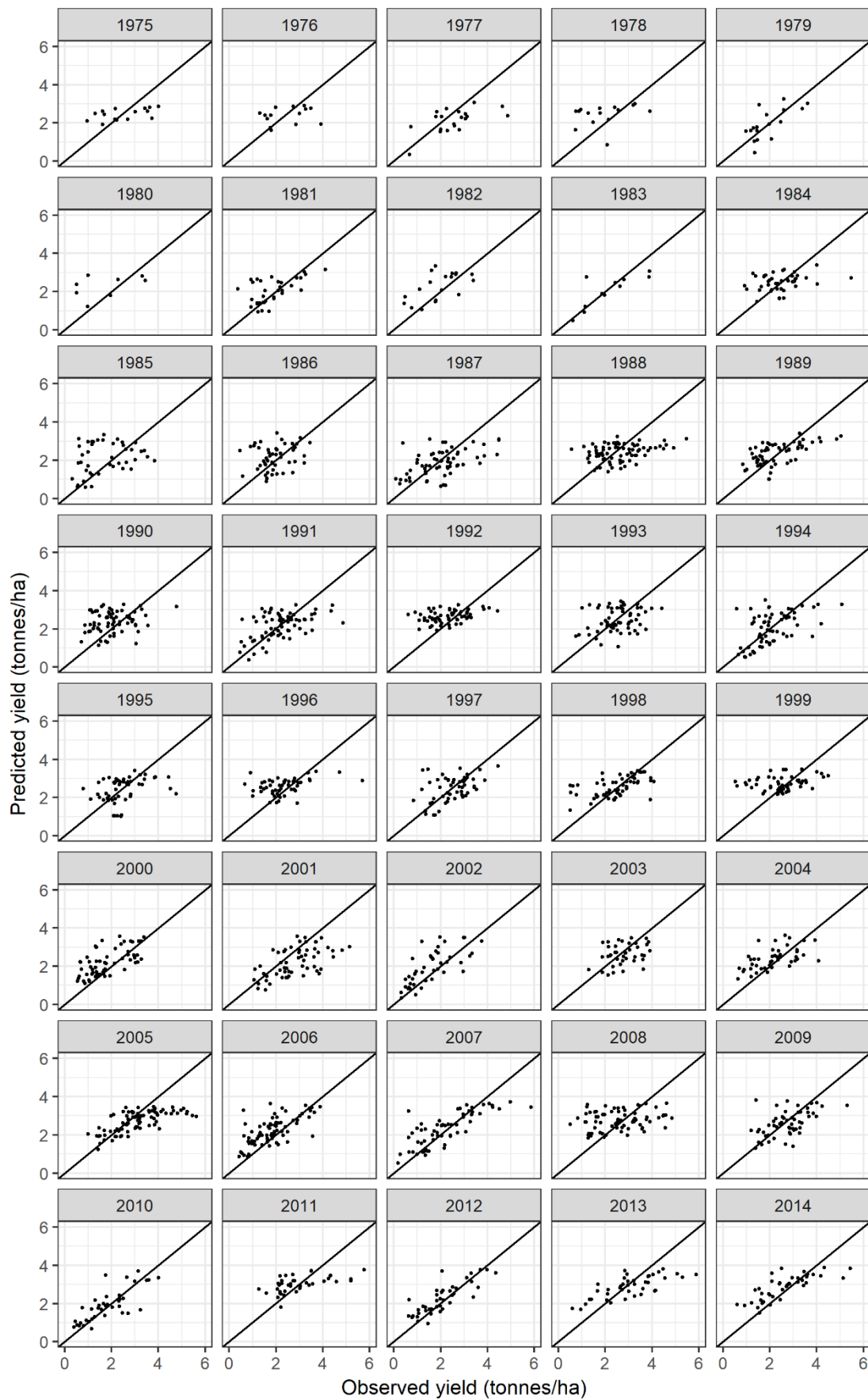
671



672

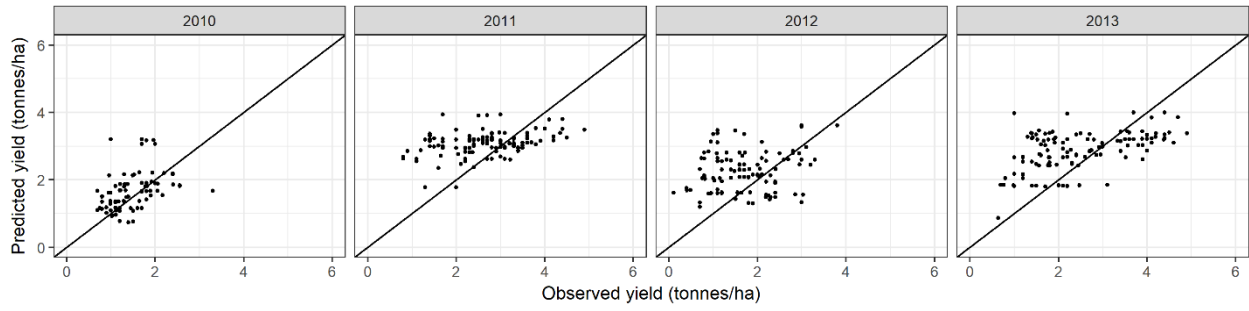
673 Figure 5: The smooth plots **a**) for $s(Wavail, k=3)$ and **b**) for $s(lat, longi, k=6)$ terms in the “best” GAM
 674 model. Note: The shaded region represents approximate pointwise 95% confidence intervals, the vertical
 675 bars at the base of the plots represent a frequency plot of the predictor variable.

676



677
 678 Figure 6: Plot of the observed versus predicted yields by year (aggregated over variety and soils). The
 679 lines show the one-to-one relationship between observed and predicted yields.

680



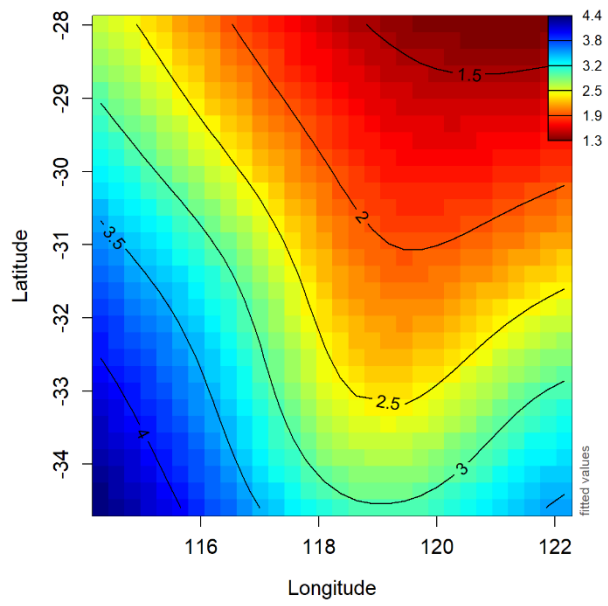
681

682 Figure 7: Plot of observed versus predicted yields by year, using the Focus Paddock data. The lines show
683 the one-to-one relationship between observed and predicted yields.

684

685

686



687

688 Figure 8: The visualised contour plot of the effect of latitude and longitude on the GAM yield prediction.

689 The contour plot was interpolated from the predicted yield in tonnes/ha, in which year is set at median of

690 the year in the variety trial data and uses the mean of *Wavail*.

692 **12. Online Supporting Information**

693 *12.1. Supplementary tables*

694 Table S1: Classification of the agricultural soil groups (Ag_soil_desc) to the functionally simplified soil
695 classes (Soil_6class).

696

697 Table S2: Description of the variables for land capability and qualities from the Natural resource
698 information.

699

700 Table S3: Mean yield and available water by year.

701

702 Table S4: List of models for water-limited yield tested with cross validation in this study.

703

704 Table S5: Results of convergence diagnostics for Gibbs simulations.

706 12.2. Supplementary figures

707 Figure S1: Map of the 164 location coordinates of Focus Paddocks with interpolation of *avgrf30*.
708

709 Figure S2: Plot of observed yield by year at the location of Newdegate Research Station (N=796) in the
710 variety trial data set. Note: The scatter plots were labelled by variety, where red hallow-circle=Calingiri,
711 green triangle=Carnamah, blue plus=Spear, gold cross=Tincurrin, purple diamond=Westonia, pink
712 inverted triangle=Wyalkatchem.
713

714 Figure S3: Model checking of the best GAM model using the residual plots including normal Q-Q plot,
715 residuals vs. fitted values, histogram of residuals and response vs. fitted values.
716

717 Figure S4: Plots showing how the correlation of annual observed (obs.yield) versus predicted yields
718 (pred.yield) varies with (a) the correlation of observed yield with *Wavail*, and (b) mean annual *Wavail*.
719 Note: the number in the circle represents the last two digits of the year.
720

721 Figure S5: The smooth plots for $s(Wavail)$ and $s(lat, longi)$ terms in the GAM model using Focus
722 Paddocks data. Note: The shaded region represents approximate pointwise 95% confidence intervals, the
723 vertical bars at the base of the plots represent a frequency plot of the predictor variable.
724

725 Figure S6: Trace and density plot for all of the parameters for the Gibbs simulations.
726

727 Figure S7: Gelman-Rubin-Brooks plot for all of the parameters for the Gibbs simulations.
728

729 Figure S8: Geweke-Brooks plot for all of the parameters for the Gibbs simulations.
730

731 Figure S9: Using the variety trial data, comparison of the predicted yields from the GAM model and the
732 potential yields from the French & Schultz approach. Where **a**) observed yields versus the French &
733 Schultz's potential yields; **b**) observed yields versus the GAM predicted yields; **c**) GAM predicted yields
734 versus French & Schultz's potential yields; **d**) water available versus French & Schultz's potential yields
735 minus the GAM predicted yields. Note: The scatter plots were labelled by soil types (red hallow-
736 circle=duplex, green triangle=clays, blue plus=gravel, gold cross=loamy, purple diamond=sandy, pink
737 inverted triangle=wet).
738

739 Figure S10: Using the Focus Paddocks data, comparison of the predicted yields from the GAM model and
740 the potential yield from the French & Schultz approach. Where **a**) observed yields versus the French &
741 Schultz's potential yields; **b**) observed yields versus the GAM predicted yields; **c**) GAM predicted yields

742 versus French & Schultz's potential yields; **d**) water available versus French & Schultz's potential yields
743 minus the GAM predicted yields. Note: The scatter plots were labelled by soil types (red hallow-
744 circle=duplex, green triangle=clays, blue plus=gravel, gold cross=loamy, purple diamond=sandy, pink
745 inverted triangle=wet).
746