

# An assessment of the predictability of column minimum dissolved oxygen concentrations in Chesapeake Bay using a machine learning model

Andrew C. Ross<sup>a,b,\*</sup>, Charles A. Stock<sup>b</sup>

<sup>a</sup>*Princeton University Program in Atmospheric and Oceanic Sciences, 300 Forrestal Road, Sayre Hall, Princeton, NJ 08540, USA*

<sup>b</sup>*NOAA Geophysical Fluid Dynamics Laboratory, Princeton University Forrestal Campus, 201 Forrestal Road, Princeton, NJ 08540, USA*

---

## Abstract

Subseasonal to seasonal forecasts have the potential to be a useful tool for managing estuarine fisheries and water quality, and with increasing skill at forecasting conditions at these time scales in the atmosphere and open ocean, skillful forecasts of estuarine salinity, temperature, and biogeochemistry may be possible. In this study, we use a machine learning model to assess the predictability of column minimum dissolved oxygen in Chesapeake Bay at a monthly time scale. Compared to previous models for dissolved oxygen and hypoxia, our model has the advantages of resolving spatial variability and fitting more flexible relationships between dissolved oxygen and the predictor variables. Using a concise set of predictors with established relationships with dissolved oxygen, we find that dissolved oxygen in a given month can be skillfully predicted with knowledge of stratification and mean temperature during the same month. Furthermore, the predictions generated by the model are consistent with expectations from prior knowledge and basic physics. The model reveals that accurate knowledge or skillful forecasts of the vertical density gradient is the key to successful prediction of dissolved oxygen, and prediction skill disappears if stratification is only known at the beginning of the forecast. The lost skill cannot be recovered by replacing stratification as a predictor with variables that have a lagged correlation with stratification (such as river discharge); however, skill is obtainable in many cases if stratification can be forecast with an error of less than about  $1 \text{ kg m}^{-3}$ . Thus, future research on hypoxia forecasting should focus on understanding and forecasting variations in stratification over subseasonal time scales

(between about two weeks and two months).

*Keywords:*

estuaries, dissolved oxygen, prediction, stratification, USA, Chesapeake Bay

---

## 1 1. Introduction

2 Chesapeake Bay, a coastal plain estuary located along the Mid-Atlantic Bight, expe-  
3 riences extensive hypoxia and anoxia in the summer following the delivery of nutrients by  
4 the spring freshet and the establishment of strong density stratification (Newcombe and  
5 Horne, 1938; Taft et al., 1980; Officer et al., 1984). Although there is some evidence that  
6 hypoxia has been an occasional feature of the bay for centuries (Karlsen et al., 2000),  
7 many studies have identified a dramatic increase in the extent and severity of hypoxia as  
8 a result of increased nutrient loading over the last century (Officer et al., 1984; Karlsen  
9 et al., 2000; Hagy et al., 2004; Murphy et al., 2011). Other estuaries and coastal systems  
10 worldwide exhibit similar increases in hypoxia, primarily as a result of increases in fer-  
11 tilizer runoff and other anthropogenic nutrient inputs (Diaz, 2001; Diaz and Rosenberg,  
12 2008; Rabalais et al., 2010; Breitburg et al., 2018). In the future, climate change and  
13 sea-level rise have the potential to alter the intensity and frequency of hypoxia, both  
14 in Chesapeake Bay (Najjar et al., 2010; Irby et al., 2018) and globally (Rabalais et al.,  
15 2010).

16 Extensive regulations have been implemented to reduce pollutants in Chesapeake  
17 Bay, including nitrogen and phosphorus, with the goal of improving water quality and  
18 reducing hypoxia (Linker et al., 2013; Shenk and Linker, 2013). Recently, there has been  
19 some evidence that water clarity and dissolved oxygen concentrations have improved  
20 (Zhang et al., 2018) and that coverage of submerged aquatic vegetation has expanded  
21 (Gurbisz and Michael Kemp, 2014; Lefcheck et al., 2018). However, historically progress  
22 has been slow (Boesch, 2006) and currently less than half of the bay area meets all water  
23 quality goals (Zhang et al., 2018).

---

\*Corresponding author

*Email addresses:* `andrew.c.ross@noaa.gov` (Andrew C. Ross), `charles.stock@noaa.gov` (Charles A. Stock)

24 While hypoxia and anoxia are nearly always present in some deep areas of Chesapeake Bay during the summer months, both the timing of hypoxia development and  
25 the spatial extent of hypoxia can vary dramatically (Hagy et al., 2004; Scully, 2016b).  
26 The susceptibility of the bay to hypoxia and the large interannual variability of hypoxia  
27 driven by weather and climate variability pose challenges for water quality and marine  
28 resource management (Boesch et al., 2001; Testa et al., 2017). Skillful forecasts of future  
29 weather and climate have the potential to improve the management of water quality  
30 and fisheries; for example, subseasonal to seasonal scale forecasts of temperature can  
31 improve the effectiveness of fisheries management (Hobday et al., 2016; Tommasi et al.,  
32 2017). Similarly, Huang and Smith (2011) show that accounting for hypoxia improves  
33 management of brown shrimp in the Neuse River Estuary; when hypoxia is more severe,  
34 the optimal opening date of the fishery is earlier in the year.  
35

36 Statistical models have been developed for forecasting the volume of hypoxic water in  
37 Chesapeake Bay (Scavia et al., 2006; Liu et al., 2011; Murphy et al., 2011), and although  
38 these forecasts are regularly published online and have received attention from the media  
39 and general public (Testa et al., 2017), the forecasts are not currently considered in man-  
40 agement of Chesapeake Bay water quality or fisheries. One key limitation is that these  
41 forecasts predict overall hypoxic volume and provide no information about the spatial  
42 distribution of hypoxia. Accounting for spatial variability is an important component  
43 of ecosystem based fisheries management (Marasco et al., 2007), and resolving spatial  
44 variability is particularly important in Chesapeake Bay because the bay straddles two  
45 states (Maryland and Virginia) and has been divided into five categories for regulation  
46 of dissolved oxygen and water quality (Batiuk et al., 2009). Additionally, although pre-  
47 vious forecast models appear to have modest skill at predicting hypoxic volume, the  
48 models have not been thoroughly evaluated for predictive skill beyond the period of data  
49 used to fit the forecast models. Therefore, the development of skillful, spatially resolved  
50 subseasonal hypoxia forecasts is an essential step for aiding and improving management  
51 decisions.

52 In this study, we assess the predictability of dissolved oxygen at a monthly time scale  
53 for many locations in Chesapeake Bay by combining a simple mechanistic set of predictors  
54 with flexible machine learning methods. Our objectives are to explore the upper bounds

55 of prediction skill (given perfect knowledge of the mechanistic drivers) and to identify  
56 key prediction bottlenecks. Previous forecasts of Chesapeake Bay hypoxia have relied on  
57 ordinary or multiple linear regression models (Murphy et al., 2011; Prasad et al., 2011;  
58 Testa et al., 2017) or on curves derived from idealized physical models (Scavia et al.,  
59 2006; Liu et al., 2011). Machine learning methods, however, have more flexibility to rep-  
60 resent nonlinearity, spatial variability, and seasonal changes in the response of dissolved  
61 oxygen to predictor variables, thus providing an opportunity for new insights. Several  
62 studies have used machine learning methods to predict hypoxia and other biogeochemical  
63 and water quality parameters in other estuaries and coastal systems. Park et al. (2015)  
64 used regression trees to estimate chlorophyll *a* given contemporaneous observations of  
65 nutrients and water temperature; they found that the regression trees were capable of  
66 representing seasonal changes in which inputs were predictive of chlorophyll concentra-  
67 tions. Thoe et al. (2014) compared the ability of a classification tree, an artificial neural  
68 network, and three regression methods to predict the presence of fecal indicator bacteria  
69 at Santa Monica Beach; they obtained the best performance with the classification tree  
70 method. Coopersmith et al. (2010) used the k-nearest neighbor (KNN) algorithm to  
71 produce one-day forecasts of hypoxia in Corpus Christi Bay. Coopersmith et al. (2010)  
72 also considered the use of regression trees, but the performance of the regression trees  
73 was worse than KNN. Tamvakis et al. (2012) found that model trees produced superior  
74 predictions of contemporaneous chlorophyll *a* compared to an artificial neural network  
75 and multiple linear regression, and Muhling et al. (2018) used model trees to predict sur-  
76 face temperature and salinity in Chesapeake Bay using projected atmospheric conditions  
77 from an ensemble of global climate models as predictors.

78 To analyze the predictability of spatially resolved dissolved oxygen in Chesapeake Bay,  
79 we use a model tree method similar to Muhling et al. (2018). As Park et al. (2015) noted  
80 for regression trees, model trees are capable of representing seasonal changes in which  
81 inputs are predictive of the response variable; this is potentially useful in Chesapeake  
82 Bay because Scully (2016b) suggested that early summer hypoxia was driven primarily  
83 by biological processes and that physical influences on hypoxia became more important  
84 later in the summer. Also, as Muhling et al. (2018) noted, model trees are capable of  
85 extrapolating outside of the range of values in the training observations (although such

86 extrapolations should be treated with caution); this is potentially useful for using the  
87 forecast model for scenario simulations to predict the effect of climate change or nutrient  
88 loading reductions on hypoxia. In Chesapeake Bay, model trees and similar methods  
89 may be more useful than time series methods, such as autoregressive models, because  
90 the inter-monthly autocorrelation of dissolved oxygen is low (Section 4.2).

91 A danger of machine learning methods is the temptation to include diverse predictors  
92 with dubious relationships to the variable being predicted. To avoid this, we focus on a  
93 distinct set of drivers that have established relationships with dissolved oxygen (Table 1).  
94 We begin by testing the predictability of dissolved oxygen under ideal conditions where  
95 we have perfect knowledge of the state of the mechanistic predictors in Table 1. Then, we  
96 reassess the skill when permutations of the predictors requiring forecasts—temperature,  
97 mean sea level and stratification—are only known at the beginning of the forecast period.  
98 This reveals stratification and, to a lesser degree, temperature, as key bottlenecks for  
99 forecasting hypoxia. We then discuss a) the accuracy of stratification forecasts required  
100 for skillful hypoxia forecasts, and b) the viability of replacing stratification as a predictor  
101 with a lagged relationship to river discharge.

## 102 **2. Methods**

103 To predict and forecast dissolved oxygen and hypoxia, we developed a machine learn-  
104 ing model that uses a model tree to predict the monthly mean, column minimum dissolved  
105 oxygen concentration (hereafter referred to as just dissolved oxygen or DO) at a given  
106 location. We refer to this model as a “mechanistic” model because the choice of pre-  
107 dictor variables in model was based on mechanisms that are known to influence DO in  
108 Chesapeake Bay. These predictor variables, the associated datasets, and the known con-  
109 nections to DO are summarized in Table 1. Based on common availability in all datasets,  
110 we used data from 1986 to 2017. These data were split into training and testing groups  
111 to fit and evaluate the model; the model was fit to the training dataset, which contained  
112 data from years 1986 to 2007, and the model was evaluated using the test dataset, which  
113 contained data for the last ten years of the record (2008 to 2017). The choice of years  
114 for training and testing does not have a substantial impact on the results; for example,  
115 using the first ten years of data as testing instead resulted in a similar model fit, and

116 although there were some differences in skill, our conclusions would not be significantly  
 117 changed.

Table 1: Variables used as inputs to the mechanistic dissolved oxygen model.

<b>Abbreviation</b>	<b>Input variable</b>	<b>Data source</b>	<b>Mechanism and references</b>
L5	TN load from Susq. River, total over previous 5 months	USGS	Phytoplankton, correlated with river discharge, estuarine circulation, and stratification.
$W_{\text{spring}}$	Mean wind along NE/SW axis, Feb-Apr	NDBC	Transport of phytoplankton biomass; Lee et al. (2013).
$\bar{T}$	Column-mean temperature anomaly, forecast month	CBP	Solubility and oxygen sinks; Li et al. (2015); Li et al. (2016).
MSL	Mean sea level anomaly, forecast month	PSMSL	Vertical exchange time, estuarine circulation, potentially correlated with stratification; Hong and Shen (2012).
$\Delta\rho$	Vertical density difference anomaly, forecast month	CBP	Mixing.
M	Forecast month		
H	Forecast hour		
D	Profile bottom depth	CBP	
X	Longitude	CBP	
Y	Latitude	CBP	

118 *2.1. Data sources and preprocessing*

119 Vertical profiles of temperature, salinity, and dissolved oxygen were obtained from the  
 120 Chesapeake Bay Program (CBP) Water Quality Database (Chesapeake Bay Program,  
 121 2018). All three variables were typically measured at 1 m intervals in each profile, and the  
 122 measurements were typically taken bimonthly for each site during the warm season. We  
 123 selected data only from sites that had frequent observations during May to September in  
 124 the last 5 years of the training period (2003 to 2007) by requiring that a site have data for

125 at least 20 of the 25 months in this time frame. We did not include sites that were located  
126 in the upper reaches of some tributaries and that never experience hypoxia (defined here  
127 as column minimum concentration below  $2 \text{ mg L}^{-1}$ ), and we also did not include a  
128 cluster of sites in the Elizabeth River near Norfolk that have experienced hypoxia in the  
129 past. We assumed that variability in dissolved oxygen in these regions is driven by more  
130 localized factors, such as discharge from minor tributaries and point source pollution,  
131 compared to the bay mainstem factors considered herein.

132 For each vertical profile, we calculated the column mean temperature and the column  
133 minimum dissolved oxygen concentration. We also obtained density from the tempera-  
134 ture and salinity profiles using the International Thermodynamic Equation Of Seawater—  
135 2010 (IOC, SCOR and IAPSO, 2010), and we calculated the density stratification as the  
136 difference between the density nearest the bottom and nearest the surface (so that a  
137 more positive value indicates a more stable density stratification).

138 We subtracted the climatological mean values from the CBP data to prevent the  
139 strong seasonal cycles of dissolved oxygen, temperature, and salinity from overwhelming  
140 the interannual variability that we seek to predict. To subtract the climatology from a  
141 variable  $y$  at a site  $i$ , we fit a generalized additive model (Hastie and Tibshirani, 1986;  
142 Wood, 2006) with a smooth seasonal cycle and a constant mean:

$$y_{ij} = s_i(DOY_j) + \beta_i + \epsilon_{ij}$$

143 where  $s_i$  is a cyclic cubic spline,  $DOY_j$  is the day of year of the  $j$ -th observation,  $\beta_i$  is the  
144 long-term mean, and  $\epsilon_{ij}$  is an independent, normally-distributed residual. A separate  
145 model was fit for each variable and site using the training dataset. The models were used  
146 to predict climatological mean values for each observation in both the training and testing  
147 datasets, and the fitted climatological values were subtracted from the observations to  
148 produce anomalies. Finally, anomalies were averaged at sites with multiple observations  
149 in a given month to produce a time series of monthly anomaly values for each site.

150 We also calculated lagged values (the value from the previous month) of the mean  
151 temperature and density stratification anomalies. At each measurement site, all data (in-  
152 cluding non-lagged variables) were eliminated if there were no measurements during the  
153 previous month. After applying this restriction and the restrictions discussed previously,  
154 126 unique locations remained in the database. A text file providing the names and

155 coordinate information of these 126 locations is provided in the supporting information.  
156 The training dataset contained 11,810 vertical profiles, and the test dataset contained  
157 4,936 profiles.

158 Data for the input of total nitrogen (TN) from the Susquehanna River were obtained  
159 from the United States Geological Survey (USGS) (Moyer and Blomquist, 2018). These  
160 data were produced by combining observations and the Weighted Regressions on Time,  
161 Discharge, and Season method (Hirsch et al., 2010). As input to the model, we used  
162 the total nitrogen loading summed over the previous five months. For a June hypoxia  
163 prediction, the previous five months are January through May, which matches the period  
164 used in other studies (Scavia et al., 2006; Liu et al., 2011; Murphy et al., 2011; Testa  
165 et al., 2017).

166 Observed wind speeds and directions were obtained from the National Data Buoy  
167 Center for Thomas Point, MD, a location in the upper Chesapeake Bay near Annapolis,  
168 MD. Winds were measured at 18 m above mean sea level. As a predictor in the models,  
169 we included mean wind speed along the northeast-southwest direction, averaged over  
170 February to April. Lee et al. (2013) suggested that winds along this axis influence the  
171 transportation of phytoplankton biomass. Because the Thomas Point station measured  
172 winds for only six days during the February to April period of 2010, the mean NE-SW  
173 wind for 2010 was determined from the value observed at Rappahannock Light, a station  
174 with similar anemometer elevation (16.9 m) located over water closer to the bay mouth.  
175 Other periods of missing data for the Thomas Point station were shorter, and the mean  
176 February-April wind was determined from all available data from the station.

177 Monthly mean sea level anomaly at Kiptopeke Beach was obtained from the Perma-  
178 nent Service for Mean Sea Level (Holgate et al., 2013). We chose this location because  
179 the data is available for the same time period as the other variables and contains less  
180 missing data than most other sites in the bay. Months that were missing in the dataset  
181 were imputed with linear interpolation.

## 182 *2.2. Model for column minimum dissolved oxygen*

183 The machine learning model for dissolved oxygen was built using a model tree (Quin-  
184 lan, 1992) as implemented and extended by the Cubist package (Kuhn et al., 2018) for  
185 R (R Core Team, 2017). In the model tree method, the training data are iteratively



186 partitioned into groups based on the values of the predictor variables, forming a tree  
187 that contains a node for each division of the data. A multiple linear regression model is  
188 developed for the data at each node of the tree, and the final predicted value is generated  
189 from a combination of the regressions along the path of the tree traversed for the given  
190 predictors (Quinlan, 1992; Kuhn et al., 2018). Model trees are controlled by a parameter  
191 for the number of “rules”, which sets the maximum number of partitions of the data  
192 included in the model. Cubist allows the addition of “neighbors” to the model, in which  
193 case the prediction for a given set of predictors is adjusted by the difference between  
194 the actual and predicted values for a specified number of neighboring, similar predictors  
195 (Quinlan, 1993). Cubist also includes the option to use “committees”, in which case  
196 the final prediction is an average of a specified number of model trees that iteratively  
197 attempt to balance errors produced by other trees (Kuhn et al., 2018).

198 We determined the approximate optimal value for each of the three parameters by  
199 searching a  $4 \times 4 \times 4$  grid containing 25, 50, 100 and 200 rules; 1, 10, 25, and 50  
200 committees; and 0, 1, 2, and 5 neighbors. Each of the 64 parameter sets was evaluated  
201 using 10-fold cross-validation, repeated 10 times, with the training dataset. The optimal  
202 set of parameters, which minimized the mean squared error of predicted DO over all  
203 stations, was 100 rules, 50 committees, and 0 neighbors.

### 204 *2.3. Model evaluation*

205 The model predictions of dissolved oxygen anomaly were compared with the obser-  
206 vations by calculating the Pearson correlation coefficient, the mean bias, and the root  
207 mean square error for each site using predictions from the test period. After clustering  
208 the sites and calculating cluster mean dissolved oxygen (Section 2.5), we also created  
209 target diagrams (Jolliff et al., 2009), which split the root mean square error (RMSE) into  
210 two components: bias, and unbiased (centered) RMSE. These components are plotted  
211 on the vertical and horizontal axes, respectively, so that the total RMSE is equivalent to  
212 the distance from the origin of the target diagram. The plots are normalized by divid-  
213 ing by the RMSE of climatological forecasts during the training period, so that a total  
214 RMSE below 1 indicates skill relative to a forecast of the training period climatology. We  
215 evaluate skill relative to climatology rather than persistence because the inter-monthly  
216 autocorrelation of dissolved oxygen is low (Section 4.2).

217 *2.4. Model sensitivity*

218 To verify that the model tree is physically reasonable and to determine the effect of  
 219 each input variable and how the model output is ultimately sensitive to the inputs, we  
 220 visualized the effects of individual terms in the model using plots of individual condi-  
 221 tional expectations (ICE) (Goldstein et al., 2015) and the average of the ICEs, known  
 222 as partial dependence (Friedman, 2001). These plots are commonly used to visualize  
 223 models where the functional form of the model is not easily interpretable. Following  
 224 Goldstein et al. (2015), the partial dependence is  $f_s = \mathbb{E}_{\mathbf{x}_c} [f(\mathbf{x}_s, \mathbf{x}_c)]$ , where  $\mathbf{x}$  is the  
 225 matrix of predictor variables,  $s$  denotes a set of one or more predictor variables for which  
 226 the partial dependence is calculated, and  $c$  is the compliment of this set (the remaining  
 227 predictor variables). In other words,  $f_s$  gives the effect of the variables in  $s$  averaged  
 228 over the other predictor variables. To calculate the partial dependence from the actual  
 229 data and model,  $f_s$  is estimated as

$$\hat{f}_s = \frac{1}{N} \sum_{i=1}^N \hat{f}(\mathbf{x}_s, \mathbf{x}_{ci})$$

230 where  $\hat{f}$  is the predicted value from the model and  $i$  denotes one of the  $N$  observations.  
 231 To reduce computational costs, we calculated  $\hat{f}_s$  for one variable at a time and for 41  
 232 evenly spaced values spanning the minimum and maximum values of  $s$  observed during  
 233 the training period. Additionally, we plotted the individual conditional expectations,  
 234 which are simply the  $N$  curves of  $\hat{f}$ . For a given plot, all curves were standardized by  
 235 subtracting the value of the curve at the minimum value of  $s$ , so that every line originates  
 236 at zero at the minimum value of  $s$ . This allows an easier comparison of the trajectories  
 237 of the ICE curves as the value of  $s$  is increased.

238 We also calculated the importance of each predictor variable in the mechanistic model.  
 239 For a given variable, the importance was determined as the percentage of the total number  
 240 of splits and regressions in the tree where the variable was used (Kuhn et al., 2018). This  
 241 provides a simple measure of how important each variable is; however, the output from a  
 242 model tree is also determined by the coefficients in each regression model along the tree,  
 243 and this is not captured by the importance metric.

244 *2.5. Station clustering*

245 To summarize the ability of the model to predict dissolved oxygen concentrations  
246 in different regions of the Bay, we grouped the CBP stations into eight clusters based  
247 on location and the percent of observations between May and September with hypoxia  
248 (or prevalence of hypoxia) (Figure 1). We first placed all stations where hypoxia never  
249 occurred into one cluster. Then, stations from the tributaries on the western side of the  
250 bay (Patuxent, Potomac, Rappahannock, and York Rivers) were assigned to clusters for  
251 their respective tributaries. Finally, stations in the mainstem (including eastern shore  
252 tributaries, which are shorter in length and have fewer stations than those on the western  
253 shore) were grouped into three clusters by applying k-means clustering to the latitude  
254 and prevalence of hypoxia over all months between May and September in the training  
255 period for each station. This neatly groups the stations into a “core hypoxic” region  
256 that experiences frequent hypoxia, an “upper bay” cluster that includes stations in the  
257 northern half of the bay that experience occasional hypoxia, and a “lower bay” cluster  
258 that includes stations in the southern half of the bay that also experience occasional  
259 hypoxia.

260 *2.6. Assessing the potential for forecasts*

261 The analyses described above assessed prediction with perfect knowledge of contem-  
262 poraneous conditions. In a forecast setting, however, the values of essential predictors  
263 are not known precisely. We thus considered three experiments to assess the potential for  
264 skillful forecasts of future dissolved oxygen concentrations. First, we assessed whether  
265 the contemporaneous variables in the mechanistic model (mean temperature anomaly,  
266 stratification anomaly, and mean sea level) can be replaced with other variables that  
267 are known in advance. We fit this “lagged” model by replacing the contemporaneous  
268 variables in the mechanistic model with the values observed during the previous month.

269 Second, from the results of the core mechanistic prediction analysis (Section 3.2),  
270 we found that accurate knowledge of stratification is the key to skillful predictions of  
271 dissolved oxygen in Chesapeake Bay. We therefore fit a “correlated” model by replacing  
272 stratification as a predictor with lagged river discharge variables that have a correla-  
273 tion with stratification. For this model, daily streamflow for the Susquehanna River at

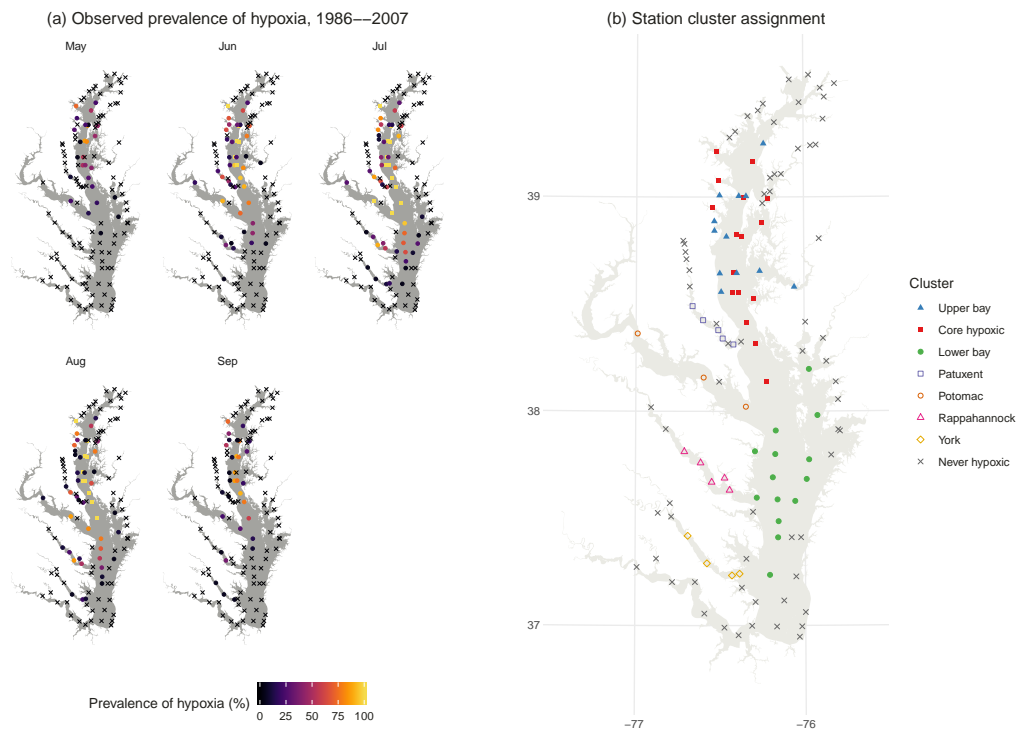


Figure 1: (a) Observed prevalence of hypoxia during the model training period. Black “x”s indicate points where hypoxia was never observed, and squares indicate points where hypoxia was always observed. Circles indicate values between these extremes. (b) Cluster assigned to each station based on geographical position and prevalence of hypoxia during May to September.

274 Conowingo, MD, the Potomac River near Washington, D.C., and the James River near  
275 Richmond, VA were obtained from the USGS. Together, these rivers represent nearly 80%  
276 of the typical freshwater discharge to the bay (Boicourt et al., 1999). The streamflow  
277 data were averaged monthly, and streamflow anomalies were calculated by subtracting  
278 the 1986 to 2007 means for each calendar month. Finally, lagged streamflow anoma-  
279 lies were calculated by taking a rolling average of the anomalies over the previous three  
280 months.

281 Lastly, to assess how accurate stratification forecasts need to be to support skillful  
282 hypoxia forecasts, we quantified the degradation of prediction skill in response to im-  
283 perfect stratification forecasts with increasing levels of noise. We ran simulations where  
284 Gaussian random noise with zero mean and various levels of variance was added to the  
285 observed stratification during the test period. The simulations assumed perfect spatial  
286 error correlation (i.e. in a given simulation, year, and month, all locations have the same  
287 error). 100 simulations were conducted for each level of error variance. For each simula-  
288 tion, we used the mechanistic model to predict dissolved oxygen using the temperature,  
289 mean sea level, spring winds, and nitrogen loading from the test dataset along with the  
290 perturbed stratification data. Then, for each region and calendar month, we calculated  
291 the average RMSE over the 100 simulations for each level of variance.

### 292 **3. Results**

#### 293 *3.1. Dissolved oxygen hindcast with mechanistic predictors*

294 With the stratification, mean temperature, and other values observed during the  
295 prediction month as inputs, the model tree produces skillful predictions of minimum  
296 dissolved oxygen anomalies during the test period. The model predictions have at least  
297 moderate correlation with the observations at the majority of sites: over all months, 54%  
298 of correlation coefficients are above 0.5 (Figure 2a). A few poor or negative correlations  
299 are found in central and lower bay along the thalweg. Except at a few stations, bias is  
300 low during the test period (Figure 2b). Over all sites and months, the predictions during  
301 the test period are essentially unbiased, with a mean bias of  $-0.07 \text{ mg L}^{-1}$  and the 25th  
302 to 75th percentiles spanning  $-0.3$  to  $0.2 \text{ mg L}^{-1}$ . The bias does tend to become more  
303 negative (i.e. model predictions are too low) as the months progress from May (mean

304 bias  $0.07 \text{ mg L}^{-1}$ ) to August (mean  $-0.2 \text{ mg L}^{-1}$ ), and several stations along the deep  
305 channel also have a large negative bias in September. Despite the biases, the overall  
306 model errors are reasonable, with predictions for 70% of all stations and months having  
307 lower RMSEs than climatological predictions (Figure 2c). RMSEs are generally low near  
308 the mouth of the bay and in some of the tributaries, with slightly higher errors present  
309 in the center of the bay. Despite low errors near the mouth of the bay, many points there  
310 are not skillful relative to climatology. This suggests that the interannual variation is low  
311 at these points, potentially as a result of exchange with saturated water from the shelf.  
312 Consistent with results from previous metrics, many points along the thalweg are also  
313 not skillful relative to climatology. In the tributaries, despite sometimes having higher  
314 RMSEs compared to average, most points are skillful relative to climatology. Overall,  
315 65% of RMSEs are below  $1 \text{ mg L}^{-1}$ , and the median RMSE is  $0.8 \text{ mg L}^{-1}$ . To put these  
316 values in context, we have included a figure of the mean minimum DO concentration for  
317 each station and month in the Supporting Information (Figure S1).

318 When aggregated to cluster means, the model predictions are generally skillful com-  
319 pared to the training period climatology (Figure 3), as indicated by points inside the  
320 solid circles. Overall, skill is highest in June through August, when all regions have  
321 lower errors than the climatological reference forecast. Most of the model predictions  
322 have lower variances than the observations (indicated by points to the left of the origin).  
323 Because the model predictions still have reasonable correlation with the observations  
324 (Figures 2a and 3), the model predictions are essentially a smoothed representation of  
325 reality. The core hypoxic cluster has lower skill than other clusters due to both larger  
326 biases than in other regions and a failure to capture the weak variability of dissolved  
327 oxygen in this region. However, because severe hypoxia is nearly always present during  
328 the summer months in this region, the lower skill would have a limited impact on pre-  
329 dicting the presence or absence of hypoxia. Predictions for the lower bay are skillful for  
330 May through August; however, skill declines significantly in September.

### 331 *3.2. Predictor importance and sensitivity*

332 On average, the model dissolved oxygen predictions are most sensitive to the vertical  
333 density stratification (Figure 4). Consistent with physical expectations, the marginal  
334 effect of increased stratification is to significantly reduce the concentration of dissolved

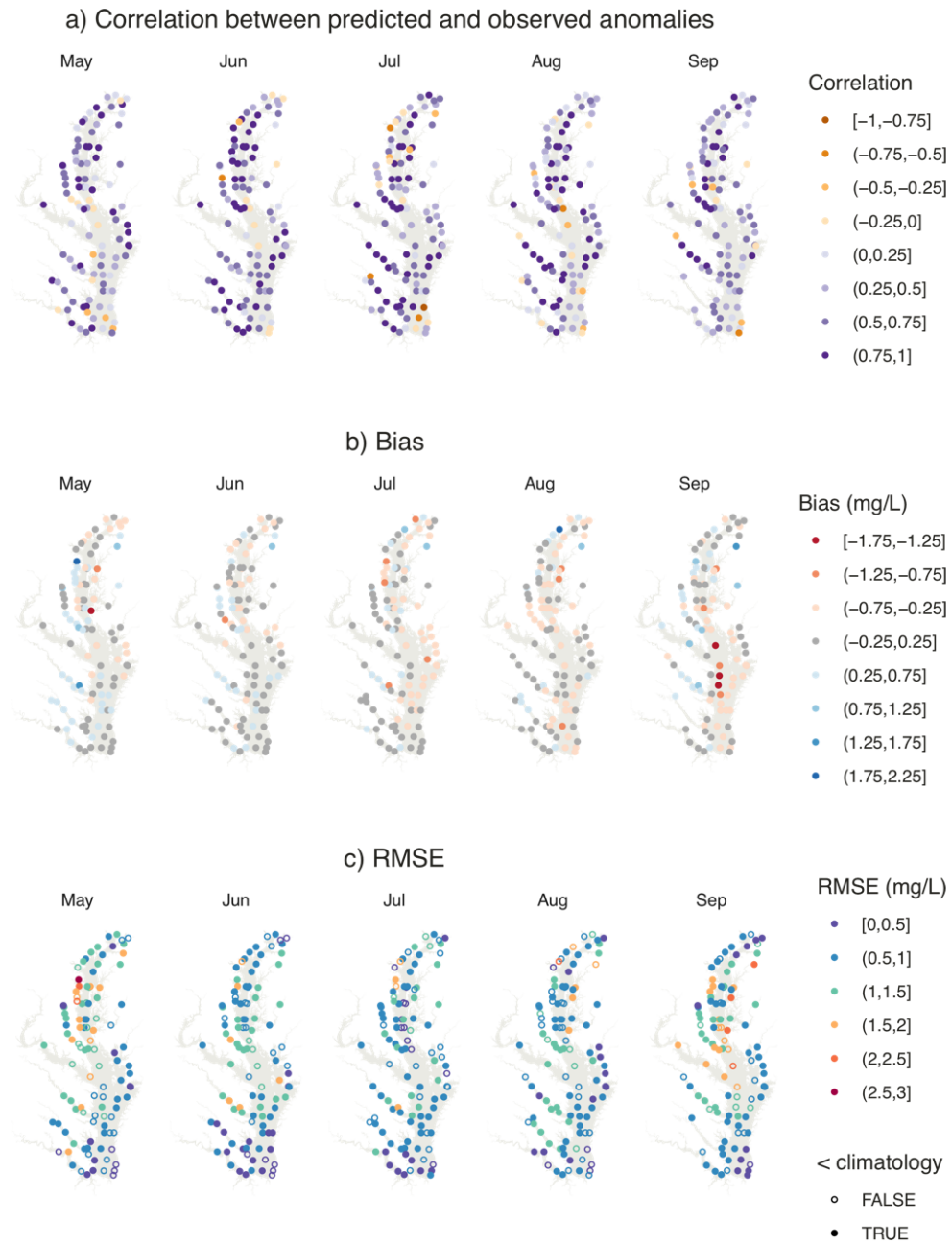


Figure 2: Skill of the main dissolved oxygen model at the station level: correlation coefficient (a), bias (b), and root mean square error (c). Solid points in panel (c) indicate lower errors than a climatological forecast.

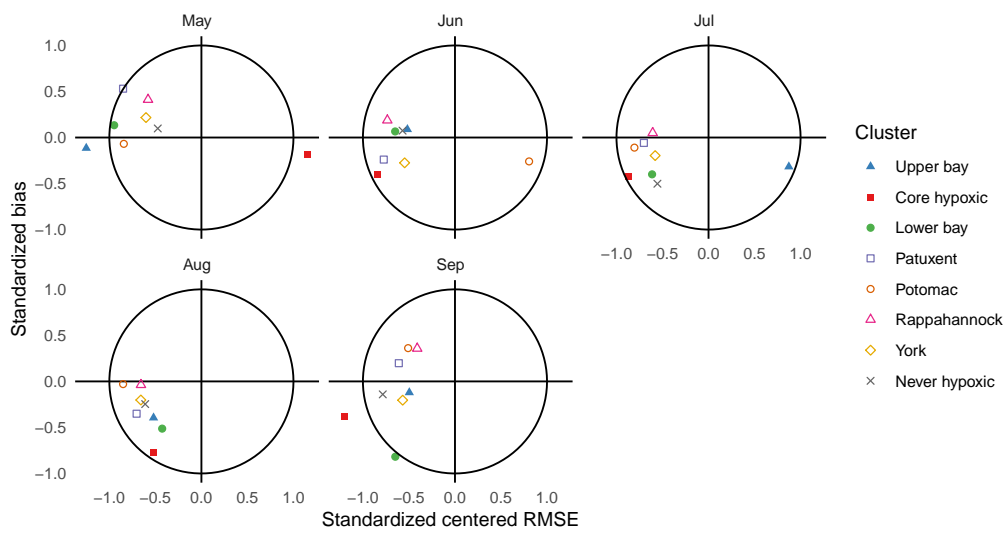


Figure 3: Target diagrams (Section 2.3) for cluster-mean predicted dissolved oxygen. Points inside the circle are considered skillful relative to the training period climatology. Points with a negative standardized centered RMSE have lower interannual variability than the observations.



335 oxygen. The ICE plots suggest that the marginal effect of stratification is stronger for  
336 some conditions or locations than others. A closer investigation showed that points where  
337 stratification has a large marginal effect in the model are typically shallow (not shown).  
338 This could be interpreted as an effect of the density gradient (an equal density difference  
339 over a shallower depth implies a higher, more stable density gradient) or a result of the  
340 lower variability of minimum dissolved oxygen in deeper regions.

341 Warmer water is modeled to have a lower dissolved oxygen concentration, which is  
342 consistent with the decreased oxygen solubility and increased biological activity associ-  
343 ated with warmer water. Unlike stratification, the effect of temperature is not a strong  
344 function of depth. The remaining variables have relatively weak effects on dissolved oxy-  
345 gen on average, although the individual conditional expectations show a fair amount of  
346 variability and suggest that interactions with other variables are present. Mean sea level  
347 and nutrient loading have weak positive effects on DO on average, while stronger winds  
348 from the northeast (positive  $W_{\text{spring}}$ ) have a weak negative effect. Although all three co-  
349 ordinate variables (depth, latitude, and longitude) have zero partial dependence because  
350 the model was fit to anomalies, the ICE plots reveal significant interactions with other  
351 variables, especially for depth and latitude. In addition to the already noted interaction  
352 between stratification and depth, interactions with latitude are not surprising: because  
353 Chesapeake Bay is roughly oriented along the north-south axis, most along-channel vari-  
354 ations, including variations in tidal amplitude and mean salinity, can be described as  
355 functions of latitude. The ICE plots for both latitude and longitude also diverge around  
356  $38.5^\circ$  and  $-76^\circ$ , respectively. This region typically has both low dissolved oxygen and  
357 frequent hypoxia along the center channel and higher dissolved oxygen and infrequent  
358 hypoxia adjacent to the channel and in the Choptank River (Figure 1). The divergence  
359 in the ICE plots suggests that the model has learned the difference between these two  
360 regions.

361 The predictor importance metric (Figure 5), which is based on the percent of the  
362 splits and regressions in the model tree in which a given variable is used, is generally  
363 consistent with the sensitivities revealed in the ICE plots. In the mechanistic model,  
364 density stratification remains the single most important variable for predicting dissolved  
365 oxygen. Latitude and depth are the two most important coordinate variables. Mean

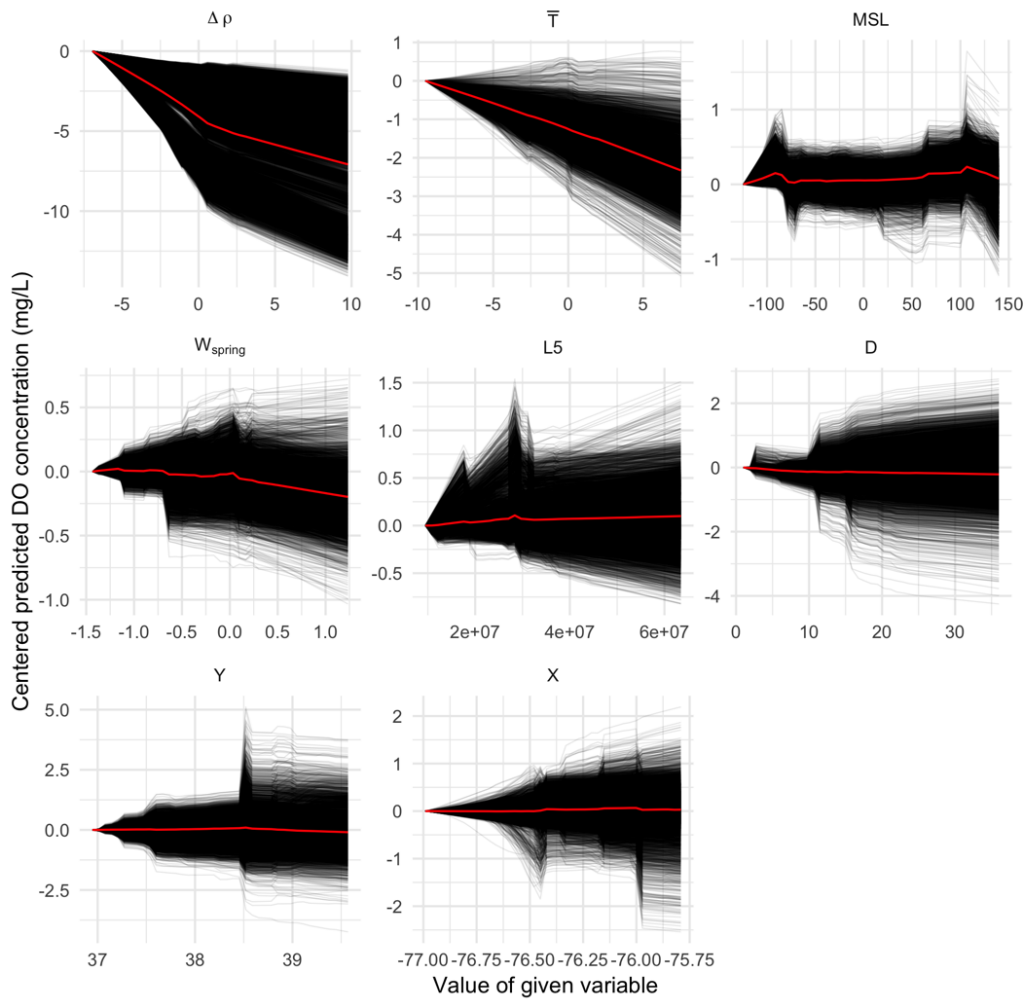


Figure 4: Individual conditional expectations (black lines) and partial dependence (red lines) for several of the predictors in the mechanistic model. Note that the y-axis for each plot is different.

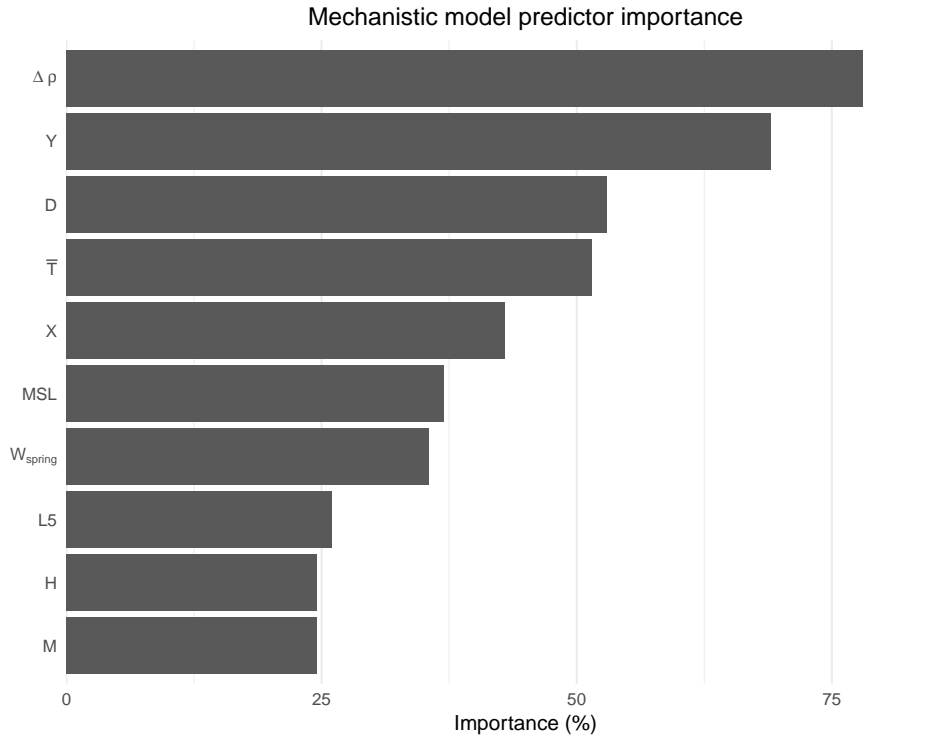


Figure 5: Importance of each variable (Section 2.4). Symbols are defined in Table 1.

366 temperature anomaly also appears in just over half of the splits and regressions, while  
 367 sea level, winds, and nitrogen loading are relatively unimportant.

### 368 3.3. Limits of predictability

369 Because the mechanistic model results show that knowledge of stratification is the  
 370 key to skillful prediction of dissolved oxygen, we consider several modifications to the  
 371 model (detailed in Section 2.6) to explore the limits of predictability of DO and to  
 372 potentially make the model useful in a forecast setting where stratification is not perfectly  
 373 predictable. First, we create a “lagged” model by replacing all contemporaneous variables  
 374 in the model (mean temperature anomaly, stratification anomaly, and mean sea level)  
 375 with the values observed during the previous month. This model has significantly reduced  
 376 skill compared to the mechanistic model (Figure 6); the predicted mean DO for all  
 377 regions has a higher error than climatology in July, and errors in the remaining months

378 are centered around climatology, with predictions in some regions having comparatively  
379 higher skill and predictions in other regions having lower skill. However, the lagged model  
380 does improve the mechanistic model prediction skill in a few cases, including in the upper  
381 bay and core hypoxic regions in May and in the core hypoxic region in September.

382     Second, in Figure 6, we test a “correlated” model by replacing the stratification  
383 predictor in the mechanistic model with discharge from three major rivers that have  
384 a lagged correlation with stratification. This model produces a modest improvement  
385 over the lagged model in many regions. The correlated model has some skill in many  
386 regions in May and September, and it improves on the mechanistic model predictions  
387 in the upper bay and core hypoxic regions in these months, suggesting there is some  
388 relationship between lagged river discharge and dissolved oxygen during the fringes of  
389 the hypoxia season. However, in nearly all regions during the main summer months, the  
390 mechanistic model performs significantly better.

391     Overall, neither the correlated model nor the lagged model appear to be viable re-  
392 placements for the mechanistic model, with the possible exception of May and September  
393 in the core hypoxic and upper bay regions. This shows that stratification is the key to  
394 successful forecasts. In Figure 7, we examine how accurately stratification must be known  
395 to allow skillful DO forecasts. Results vary by month and region, but in general the stan-  
396 dard deviation of stratification anomaly errors must be less than  $1 \text{ kg m}^{-3}$  for dissolved  
397 oxygen forecasts to be skillful in the majority of the regions (assuming the mean error is  
398 zero, i.e. the stratification forecasts are unbiased). Although seemingly small, this error  
399 is comparable to the interannual standard deviation of the stratification anomaly (Figure  
400 S2). Therefore, skillful dissolved oxygen forecasts would likely be possible if skillful fore-  
401 casts of stratification were also possible. Predictions for DO in the upper bay and never  
402 hypoxic regions are more sensitive to errors in stratification than predictions in other  
403 regions; however, these results also have lower interannual variability of stratification, so  
404 the potential for predictability remains.

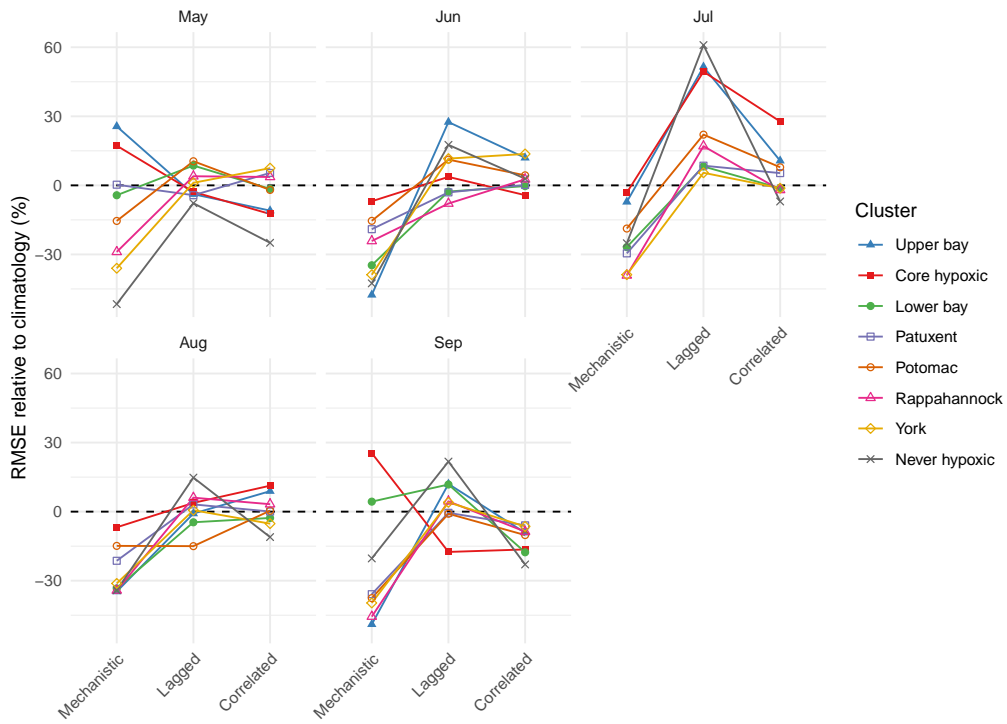


Figure 6: Root mean square error for cluster-mean dissolved oxygen during the test period. Error is normalized by the error of a prediction of climatological (training period) mean dissolved oxygen; negative values indicate errors that are lower than the climatological forecast errors. “Mechanistic” denotes predictions using the mechanistic model; “lagged” indicates predictions from a model where the contemporaneous variables in the mechanistic model are replaced with values observed in the previous month; “correlated” denotes predictions from a model similar to the mechanistic model but with the stratification anomaly replaced with correlated variables (lagged streamflow anomalies).

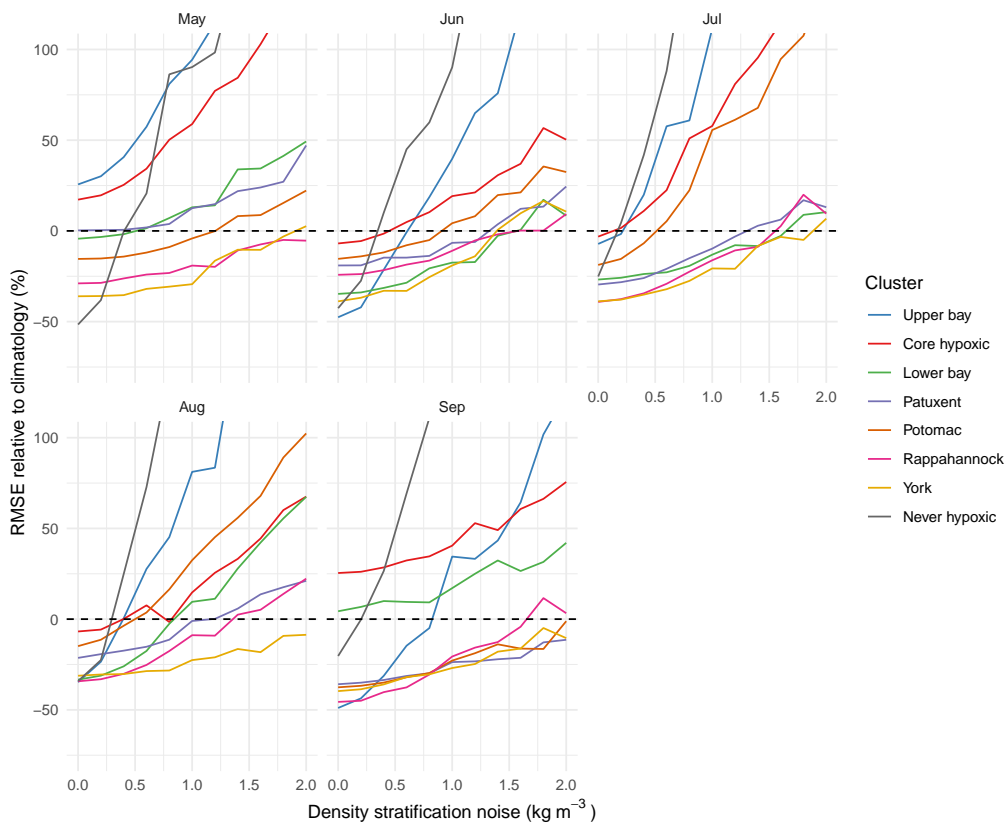


Figure 7: Root mean square error of dissolved oxygen predictions as a function of errors in the stratification anomaly input. RMSE is normalized by the error of a climatological forecast (identical to Figure 6). Stratification noise gives the standard deviation of random Gaussian errors added to the stratification predictor.

## 405 4. Discussion

### 406 4.1. Summary and comparison with previous studies

407 The mechanistic model used a concise set of predictor variables that were identified in  
408 previous studies as having a potential relationship with dissolved oxygen and hypoxia in  
409 Chesapeake Bay. Of the five time-varying variables in the model, we found that stratifi-  
410 cation and temperature had the largest influences on DO, while nutrient loading had the  
411 smallest effect. In this subsection, we summarize our findings on the effects of stratifica-  
412 tion, temperature, and nutrient loading and compare them with the results of previous  
413 studies. The comparison increases our confidence in our finding that stratification and  
414 temperature control the interannual variability of dissolved oxygen—particularly since  
415 our model, which was built on observations but with no prior assumptions about the  
416 form of the relationship between dissolved oxygen and the predictor variables, produced  
417 results that are broadly similar to other studies that have used different methods and  
418 assumptions.

#### 419 4.1.1. Stratification is the strongest predictor of dissolved oxygen

420 The mechanistic model showed that, of the variables considered, stratification is most  
421 predictive of dissolved oxygen. This is in agreement with the numerical model results in  
422 Cerco and Noel (2013); they found that stratification was the only significant predictor  
423 of bottom DO in the deeper waters of Chesapeake Bay. Our result is also partially  
424 consistent with the study of observations by Murphy et al. (2011). Murphy et al. (2011)  
425 found that stratification had a larger influence than TN load on early July hypoxic and  
426 anoxic volumes. In late July, however, Murphy et al. (2011) found that stratification  
427 had a negligible influence on hypoxia and anoxia, but stratification during the previous  
428 period (early July) had about the same influence on anoxic volumes as TN load. These  
429 findings of a strong correlation between DO and stratification are in contrast to Wang  
430 et al. (2015), who found that variability in nutrient loading was primarily responsible for  
431 interannual variability of anoxic volume. However, Wang et al. (2015) compared anoxic  
432 volume over the main bay with stratification observed at a single site (CB4.1C), whereas  
433 we have compared stratification measured at each site with concurrent dissolved oxygen  
434 measurements. Compared to Wang et al. (2015) and the other cited studies, we have also

435 considered dissolved oxygen concentrations over a broader area including the tributaries  
436 and shallow water monitoring stations.

#### 437 *4.1.2. Water temperature has a significant effect on dissolved oxygen*

438 The model in this study identified a stronger and more consistent link between warmer  
439 water and lower dissolved oxygen concentrations than previous studies have. Wang et al.  
440 (2015) found a weak negative correlation between observed summer mean bottom wa-  
441 ter temperature and anoxic volume. On the other hand, Hagy et al. (2004) found a  
442 weak positive correlation between the date of anoxia onset and the spring mean bottom  
443 temperature. Also using observed data, Scully (2016b) found essentially no correlation  
444 between summer mean sea surface temperature at Thomas Point and bay-wide hypoxic  
445 volume; however, using model simulations, Scully (2016b) found a weak positive corre-  
446 lation between temperature and hypoxic volume.

447 A possible reason that our model identified a strong and consistent link between  
448 temperature and DO is that it used column mean water temperature, which is largely  
449 independent of density stratification, as a predictor rather than using surface or bottom  
450 temperature. Modeling studies that applied long-term perturbations to atmospheric  
451 temperatures, and therefore modified the column mean temperature, have found posi-  
452 tive relationships between oxygen and temperature that are similar to this study. For  
453 example, Scully (2013) perturbed the seasonal cycle of atmospheric temperature, result-  
454 ing in a 2 °C change in water temperature and a 25% larger hypoxic volume. Irby et al.  
455 (2018) analyzed climate change simulations and concluded that the decrease in bottom  
456 DO caused by temperature change will be greater than the changes in bottom DO caused  
457 by other climate changes. Irby et al. (2018) found that the effect of temperature on solu-  
458 bility was responsible for 65-85% of the total effect of temperature on DO. Using observed  
459 data, Wang et al. (2015) also identified a weak positive correlation between atmospheric  
460 temperature and anoxic volume.

461 A second possible reason for differences between our study and some of the cited  
462 previous studies is that we included data from the tributary and shallow water regions  
463 that other studies neglected. Muller et al. (2016) found that hypoxia in two smaller  
464 tributaries, the Severn and South Rivers, was driven by temperature and temperature  
465 stratification more than by salinity and salinity stratification. However, nearly all of the



466 individual conditional expectations in Figure 4 show that increased temperature lowers  
467 DO concentration, so the effect of temperature is consistent across different stations and  
468 regions.

#### 469 *4.1.3. Nitrogen loading explains a small portion of recent oxygen variability*

470 The mechanistic model produces only a weak sensitivity of dissolved oxygen to total  
471 nitrogen loading over the study period, which is consistent with previous studies. Hagy  
472 et al. (2004) fit a linear regression to predict July hypoxic volume from January to  
473 May nitrate loading; they obtained an  $R^2$  value of 0.17. Murphy et al. (2011) fit linear  
474 regressions to predict hypoxic volume from January to May total nitrogen loading and  
475 obtained  $R^2$  values of only 0.08 and 0.21 for early and late July hypoxic volume. With  
476 only a simple model for oxygen where the oxygen consumption rate is fixed and does not  
477 respond to nutrient loading and biological activity, numerical models are still capable  
478 of skillfully simulating interannual variability in dissolved oxygen and hypoxic volume  
479 (Scully, 2010, 2013, 2016b; Irby et al., 2016). Scully (2016b) noted that despite the lack  
480 of any response to nitrogen loading in the model, the model nevertheless produced a  
481 strong correlation between nitrogen loading and hypoxic volume, which Scully (2016b)  
482 attributed to the increased stratification caused by higher discharge.

483 It is important to note that although nitrogen loading has only a weak effect on dis-  
484 solved oxygen in our model, this does not mean that efforts to reduce nitrogen loading  
485 to the bay are not worthwhile. First, of the ten predictor variables in the mechanistic  
486 model (Table 1), nitrogen loading is the only variable over which humans have some  
487 degree of control. Second, the recent interannual variability of nitrogen loading is small  
488 compared to the targeted reduction of over 40% (Cercio and Noel, 2013; Linker et al.,  
489 2013). In simple simulations using the mechanistic model with nitrogen loading uni-  
490 formly reduced by 40% over the training period, predicted dissolved oxygen increased  
491 significantly, especially over the core hypoxic region (not shown).

#### 492 *4.2. Drivers of oxygen variability not captured by the model*

493 The ability to predict dissolved oxygen using the model in this study is likely to  
494 be limited by short-term variability that is not captured in the model. Observations  
495 have shown that DO concentrations can fluctuate by several  $\text{mg L}^{-1}$  over time scales

496 as short as 5 to 15 minutes (Breitburg, 1990; Sanford et al., 1990). These fluctuations  
497 are driven by several physical factors, including barotropic tides (Breitburg, 1990) and  
498 oscillations of the pycnocline caused by internal tides and winds (Sanford et al., 1990).  
499 The short time scales associated with these events, as well as the role of advection from  
500 nearby regions, make these fluctuations essentially unpredictable using the model in this  
501 study. Because the minimum dissolved oxygen and the stratification and temperature  
502 predictors are typically derived from the average of two vertical profiles per month for  
503 each measuring site, extreme short-term variability could have also obscured the effects  
504 of the predictors in the training and testing data.

505 Modeling studies (Scully, 2010; Li and Li, 2012) and observations (Scully, 2016a)  
506 have also shown the role of winds in driving oxygen variability over time scales of a few  
507 days. Some aspects of this variability could be captured in the mechanistic model; for  
508 example, stratification also responds to these wind events (Scully et al., 2005; Li and  
509 Li, 2011; Xie and Li, 2018). However, when we constructed models that replaced the  
510 stratification predictor with various combinations of wind speed and direction averaged  
511 over the forecast month, the models did not achieve significant skill at predicting dissolved  
512 oxygen. We did not examine skill using wind predictors aggregated over shorter time  
513 scales because these winds are essentially unpredictable more than a few days in advance.

514 An additional potential source of variability and predictability that would not be  
515 captured by the model in this study is persistence of dissolved oxygen concentrations  
516 from the previous month. However, the inter-monthly correlation of dissolved oxygen  
517 in Chesapeake Bay is typically low (Figure S3). Over all months and regions, the only  
518 correlation coefficient above 0.5 is between August and September DO in the lower bay  
519 region. There is some evidence for higher correlation between months near the beginning  
520 and end of the hypoxia season (May—June and August—September). However, even  
521 in these months the correlation coefficients are typically between 0.2 and 0.4, and in  
522 other months the coefficients are even lower. Not surprisingly, using the minimum DO  
523 concentration observed during the previous month as a predictor in the model did not  
524 increase the prediction skill.

525 *4.3. Potential changes in the relationship between oxygen and predictor variables over*  
526 *time*

527 The suitability of the machine learning model for predicting future conditions could  
528 be restricted by the potential for nonstationarity in the response of oxygen to the forcing  
529 variables. For example, some estimates have found that the amount of summer hypoxia  
530 produced for a given amount of spring nitrogen loading nearly doubled during the study  
531 period (Hagy et al., 2004; Testa and Kemp, 2012). Observations also indicate that  
532 hypoxic volumes are increasing in the early summer, but volumes are decreasing in the  
533 late summer and hypoxia is breaking up earlier (Murphy et al., 2011). Given the trends  
534 in temperature, mean sea level, stratification, and other physical forcings (Murphy et al.,  
535 2011; Du et al., 2018), identifying the cause of the nonstationarity has been challenging  
536 and several hypotheses have been proposed.

537 In simulations with numerical models, a trend towards earlier development of hypoxia  
538 is consistent with the effect of warmer water (Irby et al., 2018). In this case, it would  
539 be possible to capture this effect with the mechanistic model used here. Murphy et al.  
540 (2011) suggest that an increasing trend in the strength of stratification explains some of  
541 the nonstationarity in hypoxia, which would also be captured by the mechanistic model.  
542 However, Testa and Kemp (2012) and Testa et al. (2018) proposed that these trends are  
543 a result of changes in nitrogen cycling in the bay as a result of long term hypoxia, which  
544 would not be captured by the model used in this paper.

545 In the mechanistic model, biases became negative during the test period from May  
546 to August, especially in the core hypoxic and lower bay regions (Section 3.1). This is  
547 consistent with hypoxia breaking up earlier in the test period than during the training  
548 period, and suggests that the causes of the earlier breakup are not captured by the  
549 predictors included in the mechanistic model. Despite this potential nonstationarity, the  
550 model predictions were still skillful compared to climatology during the test period, which  
551 suggests that potential nonstationarity will not have a severe impact on model predictions  
552 for the near future. Furthermore, as additional observations are collected, the model  
553 can be adapted to any nonstationarity by including these observations and adding any  
554 variables that are discovered to be causing changes in dissolved oxygen concentrations.

555 *4.4. Comparison of machine learning and other models*

556 While it was not our objective to conduct a comprehensive intercomparison of different  
557 methods for modeling dissolved oxygen, in this section we briefly discuss what our work  
558 shows may be advantages and disadvantages of the modeling approach used in this study  
559 compared to both simpler linear regression models and more complex numerical models.

560 Compared to simpler linear regression models, model trees and other machine learn-  
561 ing methods have a number of potential advantages. For example, the model trees used  
562 in this study were able to model dissolved oxygen in different months by including the  
563 calendar month as a predictor variable, which was used by the model tree algorithm as a  
564 criterion for dividing the data and fitting different regressions. Some studies using linear  
565 regression models have adopted a similar, but manual, approach by creating multiple  
566 models for different months (e.g., Testa et al. (2017)). Unlike linear regression models,  
567 model trees and many other methods are capable of fitting complex and nonlinear rela-  
568 tionships between the predictors and the variable being predicted. These advantages can  
569 lead to improved prediction skill over linear regression; for example, when we ran simple  
570 experiments using a multiple linear regression model with the same predictors as the  
571 mechanistic model, the linear regression model had lower skill in the majority of cases.  
572 However, complex machine learning models do have disadvantages compared to linear  
573 regression. The complex models can be much less interpretable, and the larger number  
574 of parameters in the complex models requires the availability of more data for training.

575 Although machine learning models can be complex, they still have advantages over  
576 even more complex numerical biogeochemical models. One clear advantage is computa-  
577 tional cost: once optimal parameters have been found using cross-validation (which takes  
578 a few hours on a quad core computer), the model tree used in this study can be trained  
579 and used to predict years of data in a few seconds. By comparison, we have used a 3D  
580 numerical model of Chesapeake Bay in other research that requires over an hour to sim-  
581 ulate a single month using a similar computer. A second advantage is the fewer number  
582 of parameters and the simpler process for learning these parameters. One disadvantage  
583 is that numerical models, which are rooted in fundamental physical principles, are more  
584 reliable when extrapolating beyond the range of historically observed conditions (for ex-  
585 ample, when simulating the effects of climate change). Numerical models also provide

586 predictions of multiple variables simultaneously and allow an easier understanding of the  
587 physical reasoning behind the predictions. Overall, the mechanistic model tree appears  
588 to have skill that is comparable to the skill that Irby et al. (2016) obtained in a compar-  
589 ison of hindcast simulations from coupled numerical biogeochemical models, although a  
590 more detailed comparison is needed.

## 591 **5. Conclusions**

592 We developed a machine learning model to forecast and predict spatially explicit min-  
593 imum dissolved oxygen in Chesapeake Bay at monthly time scales. The model results  
594 show that accurate knowledge of density stratification is the key to skillful predictions of  
595 dissolved oxygen. We developed two alternative models that replaced density stratifica-  
596 tion with other predictor variables, and neither alternative model was skillful enough to  
597 be a viable replacement for the mechanistic model. This suggests that although the mech-  
598 anistic model is capable of skillfully at predicting dissolved oxygen, accurate forecasts  
599 of stratification are necessary to use the mechanistic model to forecast future dissolved  
600 oxygen.

601 Even if machine learning models like the one used in this study are not capable  
602 of standing alone as forecast models, they have a number of potential uses, including  
603 serving as replacements for complex and expensive biogeochemical model components in  
604 a numerical ocean model capable of predicting stratification. With significantly reduced  
605 computational costs, additional numerical model ensembles can be run, which will likely  
606 increase the accuracy of both subseasonal forecasts and decadal scale climate simulations.

## 607 **6. Acknowledgements**

608 The authors thank Barbara Muhling for providing useful information about the model  
609 tree method and Xiao Liu and Fernando González Taboada for providing helpful reviews  
610 of this manuscript. We also thank three anonymous reviewers for providing reviews  
611 that improved this manuscript. This report was prepared by the authors under award  
612 NA14OAR4320106 from the National Oceanic and Atmospheric Administration, U.S.  
613 Department of Commerce and with funding from the NOAA Integrated Ecosystem As-  
614 sessment program. The statements, findings, conclusions, and recommendations are

615 those of the authors and do not necessarily reflect the views of the National Oceanic  
616 and Atmospheric Administration, or the U.S. Department of Commerce. Declarations  
617 of interest: none.

618 **References**

- 619 Batiuk, R.A., Breitburg, D.L., Diaz, R.J., Cronin, T.M., Secor, D.H., Thursby, G., 2009. Derivation  
620 of habitat-specific dissolved oxygen criteria for Chesapeake Bay and its tidal tributaries. *Journal of*  
621 *Experimental Marine Biology and Ecology* 381, S204–S215. doi:10.1016/j.jembe.2009.07.023.
- 622 Boesch, D., Brinsfield, R.B., Magnien, R.E., 2001. Chesapeake Bay eutrophication. *Journal of Environ-*  
623 *mental Quality* 30, 303–320.
- 624 Boesch, D.F., 2006. Scientific requirements for ecosystem-based management in the restoration of Ches-  
625 apeake Bay and Coastal Louisiana. *Ecological Engineering* 26, 6–26. doi:10.1016/j.ecoleng.2005.09.  
626 004.
- 627 Boicourt, W.C., Kuzmić, M., Hopkins, T.S., 1999. The inland sea: Circulation of Chesapeake Bay  
628 and the Northern Adriatic, in: Malone, T.C., Malej, A., Harding, L.W., Smolaka, N., Turner, R.E.  
629 (Eds.), *Ecosystems at the Land-Sea Margin: Drainage Basin to Coastal Sea*. American Geophysical  
630 Union, Washington, D. C., pp. 81–129. doi:10.1029/ce055p0081.
- 631 Breitburg, D., Levin, L.A., Oschlies, A., Grégoire, M., Chavez, F.P., Conley, D.J., Garçon, V., Gilbert,  
632 D., Gutiérrez, D., Isensee, K., Jacinto, G.S., Limburg, K.E., Montes, I., Naqvi, S.W., Pitcher, G.C.,  
633 Rabalais, N.N., Roman, M.R., Rose, K.A., Seibel, B.A., Telszewski, M., Yasuhara, M., Zhang, J., 2018.  
634 Declining oxygen in the global ocean and coastal waters. *Science* 359. doi:10.1126/science.aam7240.
- 635 Breitburg, D.L., 1990. Near-shore hypoxia in the Chesapeake Bay: Patterns and relationships among  
636 physical factors. *Estuarine, Coastal and Shelf Science* 30, 593–609. doi:10.1016/0272-7714(90)  
637 90095-9.
- 638 Cerco, C.F., Noel, M.R., 2013. Twenty-one-year simulation of Chesapeake Bay water quality using the  
639 CE-QUAL-ICM eutrophication model. *Journal of the American Water Resources Association* 49,  
640 1119–1133. doi:10.1111/jawr.12107.
- 641 Chesapeake Bay Program, 2018. CBP Water Quality Database (1984-present).  
642 [https://www.chesapeakebay.net/what/downloads/cbp\\_water\\_quality\\_database\\_1984\\_present](https://www.chesapeakebay.net/what/downloads/cbp_water_quality_database_1984_present).
- 643 Coopersmith, E.J., Minsker, B., Montagna, P., 2010. Understanding and forecasting hypoxia using  
644 machine learning algorithms. *Journal of Hydroinformatics* 13, 64. doi:10.2166/hydro.2010.015.
- 645 Diaz, R.J., 2001. Overview of Hypoxia around the World. *Journal of Environment Quality* 30, 275.  
646 doi:10.2134/jeq2001.302275x.
- 647 Diaz, R.J., Rosenberg, R., 2008. Spreading Dead Zones and Consequences for Marine Ecosystems.  
648 *Science* 321, 926–929. doi:10.1126/science.1156401.
- 649 Du, J., Shen, J., Park, K., Wang, Y.P., Yu, X., 2018. Worsened physical condition due to climate change  
650 contributes to the increasing hypoxia in Chesapeake Bay. *Science of the Total Environment* 630,  
651 707–717. doi:10.1016/j.scitotenv.2018.02.265.
- 652 Friedman, J.H., 2001. Greedy function approximation: a gradient boosting machine. *The Annals of*  
653 *Statistics* 29, 1189–1232.
- 654 Goldstein, A., Kapelner, A., Bleich, J., Pitkin, E., 2015. Peeking Inside the Black Box: Visualizing  
655 Statistical Learning With Plots of Individual Conditional Expectation. *Journal of Computational*  
656 *and Graphical Statistics* 24, 44–65. doi:10.1080/10618600.2014.907095.

657 Gurbisz, C., Michael Kemp, W., 2014. Unexpected resurgence of a large submersed plant bed  
658 in Chesapeake Bay: Analysis of time series data. *Limnology and Oceanography* 59, 482–494.  
659 doi:10.4319/lo.2014.59.2.0482.

660 Hagy, J.D., Boynton, W.R., Keefe, C.W., Wood, K.V., 2004. Hypoxia in Chesapeake Bay, 1950-2001:  
661 Long-term change in relation to nutrient loading and river flow. *Estuaries* 4, 634–658. doi:10.1007/  
662 BF02907650.

663 Hastie, T., Tibshirani, R., 1986. Generalized Additive Models. *Statistical Science* 1, 297–318. doi:10.  
664 1214/ss/1177013604.

665 Hirsch, R.M., Moyer, D.L., Archfield, S.A., 2010. Weighted regressions on time, discharge, and season  
666 (WRTDS), with an application to Chesapeake Bay river inputs. *Journal of the American Water*  
667 *Resources Association* 46, 857–880. doi:10.1111/j.1752-1688.2010.00482.x.

668 Hobday, A.J., Spillman, C.M., Paige Eveson, J., Hartog, J.R., 2016. Seasonal forecasting for decision  
669 support in marine fisheries and aquaculture. *Fisheries Oceanography* 25, 45–56. doi:10.1111/fog.  
670 12083.

671 Holgate, S.J., Matthews, A., Woodworth, P.L., Rickards, L.J., Tamisiea, M.E., Bradshaw, E., Foden,  
672 P.R., Gordon, K.M., Jevrejeva, S., Pugh, J., 2013. New Data Systems and Products at the  
673 Permanent Service for Mean Sea Level. *Journal of Coastal Research* 288, 493–504. doi:10.2112/  
674 JCOASTRES-D-12-00175.1.

675 Hong, B., Shen, J., 2012. Responses of estuarine salinity and transport processes to potential future  
676 sea-level rise in the Chesapeake Bay. *Estuarine, Coastal and Shelf Science* 104-105, 33–45. doi:10.  
677 1016/j.ecss.2012.03.014.

678 Huang, L., Smith, M.D., 2011. Management of an annual fishery in the presence of ecological stress: The  
679 case of shrimp and hypoxia. *Ecological Economics* 70, 688–697. doi:10.1016/j.ecolecon.2010.11.003.

680 IOC, SCOR and IAPSO, 2010. The international thermodynamic equation of seawater—2010: Calcula-  
681 tion and use of thermodynamic properties. Technical Report Manuals and Guides No. 56. Intergov-  
682 ernmental Oceanographic Commission.

683 Irby, I., Friedrichs, M.A., Da, F., Hinson, K., 2018. The competing impacts of climate change and  
684 nutrient reductions on dissolved oxygen in Chesapeake Bay. *Biogeosciences* 15, 2649–2668. doi:10.  
685 5194/bg-15-2649-2018.

686 Irby, I.D., Friedrichs, M.A., Friedrichs, C.T., Bever, A.J., Hood, R.R., Lanerolle, L.W., Li, M., Linker, L.,  
687 Scully, M.E., Sellner, K., Shen, J., Testa, J., Wang, H., Wang, P., Xia, M., 2016. Challenges associated  
688 with modeling low-oxygen waters in Chesapeake Bay: A multiple model comparison. *Biogeosciences*  
689 13. doi:10.5194/bg-13-2011-2016.

690 Jolliff, J.K., Kindle, J.C., Shulman, I., Penta, B., Friedrichs, M.A.M., Helber, R., Arnone, R.A., 2009.  
691 Summary diagrams for coupled hydrodynamic-ecosystem model skill assessment. *Journal of Marine*  
692 *Systems* 76, 64–82. doi:10.1016/j.jmarsys.2008.05.014.

693 Karlson, A.W., Cronin, T.M., Ishman, S.E., Willard, D.A., Kerhin, R., Holmes, C.W., Marot, M., 2000.  
694 Historical Trends in Chesapeake Bay Dissolved Oxygen Based on Benthic Foraminifera from Sediment  
695 Cores. *Estuaries* 23, 488. doi:10.2307/1353141.



- 696 Kuhn, M., Weston, S., Keefer, C., Coulter, N., Quinlan, R., 2018. Cubist: Rule- And Instance-Based  
697 Regression Modeling, R package version 0.2.1. <https://topepo.github.io/Cubist>.
- 698 Lee, Y.J., Boynton, W.R., Li, M., Li, Y., 2013. Role of Late Winter–Spring Wind Influencing Summer  
699 Hypoxia in Chesapeake Bay. *Estuaries and Coasts* 36, 683–696. doi:10.1007/s12237-013-9592-5.
- 700 Lefcheck, J.S., Orth, R.J., Dennison, W.C., Wilcox, D.J., Murphy, R.R., Keisman, J., Gurbisz, C.,  
701 Hannam, M., Landry, J.B., Moore, K.A., Patrick, C.J., Testa, J., Weller, D.E., Batiuk, R.A., 2018.  
702 Long-term nutrient reductions lead to the unprecedented recovery of a temperate coastal region.  
703 *Proceedings of the National Academy of Sciences* , 201715798doi:10.1073/pnas.1715798115.
- 704 Li, M., Lee, Y.J., Testa, J.M., Li, Y., Ni, W., Kemp, W.M., Di Toro, D.M., 2016. What drives interannual  
705 variability of hypoxia in Chesapeake Bay: Climate forcing versus nutrient loading? *Geophysical*  
706 *Research Letters* 43, 2127–2134. doi:10.1002/2015GL067334.
- 707 Li, Y., Li, M., 2011. Effects of winds on stratification and circulation in a partially mixed estuary.  
708 *Journal of Geophysical Research: Oceans* 116. doi:10.1029/2010JC006893.
- 709 Li, Y., Li, M., 2012. Wind-driven lateral circulation in a stratified estuary and its effects on the along-  
710 channel flow. *Journal of Geophysical Research: Oceans* 117. doi:10.1029/2011JC007829.
- 711 Li, Y., Li, M., Kemp, W.M., 2015. A Budget Analysis of Bottom-Water Dissolved Oxygen in Chesapeake  
712 Bay. *Estuaries and Coasts* 38, 2132–2148. doi:10.1007/s12237-014-9928-9.
- 713 Linker, L.C., Batiuk, R.A., Shenk, G.W., Cerco, C.F., 2013. Development of the Chesapeake Bay  
714 watershed total maximum daily load allocation. *Journal of the American Water Resources Association*  
715 49, 986–1006. doi:10.1111/jawr.12105.
- 716 Liu, Y., Arhonditsis, G.B., Stow, C.a., Scavia, D., 2011. Predicting the Hypoxic-Volume in Chesapeake  
717 Bay With the Streeter–Phelps Model : a Bayesian Approach. *Journal Of The American Water*  
718 *Resources Association* 1, 1348–1363.
- 719 Marasco, R.J., Goodman, D., Grimes, C.B., Lawson, P.W., Punt, A.E., Quinn II, T.J., 2007. Ecosystem-  
720 based fisheries management: some practical suggestions. *Canadian Journal of Fisheries and Aquatic*  
721 *Sciences* 64, 928–939. doi:10.1139/f07-062.
- 722 Moyer, D.L., Blomquist, J.D., 2018. Nitrogen, phosphorus, and suspended-sediment loads and trends  
723 measured at the Chesapeake Bay River Input Monitoring stations: Water years 1985-2017. doi:10.  
724 5066/P96NUK3Q.
- 725 Muhling, B.A., Gaitán, C.F., Stock, C.A., Saba, V.S., Tommasi, D., Dixon, K.W., 2018. Potential  
726 Salinity and Temperature Futures for the Chesapeake Bay Using a Statistical Downscaling Spatial  
727 Disaggregation Framework. *Estuaries and Coasts* 41, 349—372. doi:10.1007/s12237-017-0280-8.
- 728 Muller, A.C., Muller, D.L., Muller, A., 2016. Resolving spatiotemporal characteristics of the seasonal  
729 hypoxia cycle in shallow estuarine environments of the Severn River and South River, MD, Chesapeake  
730 Bay, USA. *Heliyon* 2, e00157. doi:10.1016/j.heliyon.2016.e00157.
- 731 Murphy, R.R., Kemp, W.M., Ball, W.P., 2011. Long-Term Trends in Chesapeake Bay Seasonal Hy-  
732 poxia, Stratification, and Nutrient Loading. *Estuaries and Coasts* 34, 1293–1309. doi:10.1007/  
733 s12237-011-9413-7.
- 734 Najjar, R.G., Pyke, C.R., Adams, M.B., Breitburg, D., Hershner, C., Kemp, M., Howarth, R., Mul-

735 holland, M.R., Paolisso, M., Secor, D., Sellner, K., Wardrop, D., Wood, R., 2010. Potential  
736 climate-change impacts on the Chesapeake Bay. *Estuarine, Coastal and Shelf Science* 86, 1–20.  
737 doi:10.1016/j.ecss.2009.09.026.

738 Newcombe, C.L., Horne, W.A., 1938. Oxygen-poor waters of the Chesapeake Bay. *Science* 88, 80–81.  
739 doi:10.1126/science.88.2273.80.

740 Officer, C.B., Biggs, R.B., Taft, J.L., Cronin, L.E., Tyler, M.A., Boynton, W.R., 1984. Chesapeake Bay  
741 Anoxia: Origin, Development, and Significance. *Science* 223, 22–27. doi:10.1126/science.223.4631.  
742 22.

743 Park, Y., Pachepsky, Y.A., Cho, K.H., Jeon, D.J., Kim, J.H., 2015. Stressor-response modeling using  
744 the 2D water quality model and regression trees to predict chlorophyll-a in a reservoir system. *Journal*  
745 *of Hydrology* 529, 805–815. doi:10.1016/j.jhydrol.2015.09.002.

746 Prasad, M.B.K., Long, W., Zhang, X., Wood, R.J., Murtugudde, R., 2011. Predicting dissolved oxygen  
747 in the Chesapeake Bay: Applications and implications. *Aquatic Sciences* 73, 437–451. doi:10.1007/  
748 s00027-011-0191-x.

749 Quinlan, J.R., 1992. Learning with continuous classes. *Machine Learning* 92, 343–348. doi:10.1.1.34.  
750 885.

751 Quinlan, J.R., 1993. Combining Instance-Based and Model-Based Learning. *Machine Learning* 76,  
752 236–243.

753 R Core Team, 2017. R: A Language and Environment for Statistical Computing. R Foundation for  
754 Statistical Computing. Vienna, Austria.

755 Rabalais, N.N., Díaz, R.J., Levin, L.A., Turner, R.E., Gilbert, D., Zhang, J., 2010. Dynamics and distri-  
756 bution of natural and human-caused hypoxia. *Biogeosciences* 7, 585–619. doi:10.5194/bg-7-585-2010.

757 Sanford, L.P., Sellner, K.G., Breitburg, D.L., 1990. Covariability of dissolved oxygen with physical  
758 processes in the summertime Chesapeake Bay. *Journal of Marine Research* 48, 567–590. doi:10.1357/  
759 002224090784984713.

760 Scavia, D., Kelly, E.L.A., Hagy, J.D., 2006. A simple model for forecasting the effects of nitrogen loads  
761 on Chesapeake Bay hypoxia. *Estuaries and Coasts* doi:10.1007/BF02784292.

762 Scully, M.E., 2010. Wind Modulation of Dissolved Oxygen in Chesapeake Bay. *Estuaries and Coasts*  
763 33, 1164–1175. doi:10.1007/s12237-010-9319-9.

764 Scully, M.E., 2013. Physical controls on hypoxia in Chesapeake Bay: A numerical modeling study.  
765 *Journal of Geophysical Research: Oceans* 118, 1239–1256. doi:10.1002/jgrc.20138.

766 Scully, M.E., 2016a. Mixing of dissolved oxygen in Chesapeake Bay driven by the interaction between  
767 wind-driven circulation and estuarine bathymetry. *Journal of Geophysical Research: Oceans* 121,  
768 5639–5654. doi:10.1002/2016JC011924.

769 Scully, M.E., 2016b. The contribution of physical processes to inter-annual variations of hypoxia in  
770 Chesapeake Bay: A 30-yr modeling study. *Limnology and Oceanography* 61, 2243–2260. doi:10.  
771 1002/lno.10372.

772 Scully, M.E., Friedrichs, C., Brubaker, J., 2005. Control of estuarine stratification and mixing by wind-  
773 induced straining of the estuarine density field. *Estuaries* 28, 321–326. doi:10.1007/BF02693915.

- 774 Shenk, G.W., Linker, L.C., 2013. Development and application of the 2010 Chesapeake Bay Watershed  
775 total maximum daily load model. *Journal of the American Water Resources Association* 49, 1042–  
776 1056. doi:10.1111/jawr.12109.
- 777 Taft, J.L., Taylor, W.R., Hartwig, E.O., Loftus, R., 1980. Seasonal Oxygen Depletion in Chesapeake  
778 Bay. *Estuaries* 3, 242. doi:10.2307/1352079.
- 779 Tamvakis, A., Miritzis, J., Tsirtsis, G., Spyropoulou, A., Spatharis, S., 2012. Effects of meteorological  
780 forcing on coastal eutrophication: Modeling with model trees. *Estuarine, Coastal and Shelf Science*  
781 115, 210–217. doi:10.1016/J.ECSS.2012.09.003.
- 782 Testa, J.M., Clark, J.B., Dennison, W.C., Donovan, E.C., Fisher, A.W., Ni, W., Parker, M., Scavia,  
783 D., Spitzer, S.E., Waldrop, A.M., Vargas, V.M.D., Ziegler, G., 2017. Ecological Forecasting and the  
784 Science of Hypoxia in Chesapeake Bay. *BioScience* 67, 614–626. doi:10.1093/biosci/bix048.
- 785 Testa, J.M., Kemp, W.M., 2012. Hypoxia-induced shifts in nitrogen and phosphorus cycling in Ches-  
786 peake Bay. *Limnology and Oceanography* 57, 835–850. doi:10.4319/lo.2012.57.3.0835.
- 787 Testa, J.M., Kemp, W.M., Boynton, W.R., 2018. Season-specific trends and linkages of nitrogen and  
788 oxygen cycles in Chesapeake Bay. *Limnology and Oceanography* doi:10.1002/lno.10823.
- 789 Thoe, W., Gold, M., Griesbach, A., Grimmer, M., Taggart, M., Boehm, A., 2014. Predicting water  
790 quality at Santa Monica Beach: Evaluation of five different models for public notification of unsafe  
791 swimming conditions. *Water Research* 67, 105–117. doi:10.1016/J.WATRES.2014.09.001.
- 792 Tommasi, D., Stock, C.A., Pegion, K., Vecchi, G.A., Methot, R.D., Alexander, M.A., Checkley, D.M.,  
793 2017. Improved management of small pelagic fisheries through seasonal climate prediction. *Ecological*  
794 *Applications* doi:10.1002/eap.1458.
- 795 Wang, P., Wang, H., Linker, L., 2015. Relative Importance of Nutrient Load and Wind on Regulating  
796 Interannual Summer Hypoxia in the Chesapeake Bay. *Estuaries and Coasts* 38, 1048–1061. doi:10.  
797 1007/s12237-014-9867-5.
- 798 Wood, S.N., 2006. *Generalized additive models: an introduction with R*. Chapman and Hall/CRC, New  
799 York.
- 800 Xie, X., Li, M., 2018. Effects of Wind Straining on Estuarine Stratification: A Combined Ob-  
801 servational and Modeling Study. *Journal of Geophysical Research: Oceans* 123, 2363–2380.  
802 doi:10.1002/2017JC013470.
- 803 Zhang, Q., Murphy, R.R., Tian, R., Forsyth, M.K., Trentacoste, E.M., Keisman, J., Tango, P.J., 2018.  
804 Chesapeake Bay's water quality condition has been recovering: Insights from a multimetric indicator  
805 assessment of thirty years of tidal monitoring data. *Science of the Total Environment* 637–638, 1617–  
806 1625. doi:10.1016/j.scitotenv.2018.05.025.

**Training observations**



**Mechanistic model**

Temperature, stratification, wind,  
mean sea level, nitrogen loading

**Predictor importance  
and sensitivity**

**Lagged model**

Replace temperature, stratification,  
MSL with values from previous month

**Correlated model**

Replace stratification  
with prior river discharge

**Testing observations**



**Skill assessment**

