

1 **A Bayesian Approach to Recreational Water Quality**
2 **Model Validation and Comparison in the Presence of**
3 **Measurement Error**

4 **E. Potash^{1*} and S. Steinschneider²**

5 ¹University of Illinois Urbana-Champaign, Illinois, USA
6 ²Cornell University, Ithaca, New York, USA

7 **Key Points:**

- 8 • Recreational water quality models are typically validated and compared ignoring
9 measurement errors
10 • New methods to account for such errors are developed and applied to fecal
11 indicator bacteria models at Chicago beaches
12 • Accounting for measurement error significantly affects validation results and
13 reveals high uncertainty

*1101 W. Peabody, Suite 350 (NSRC), MC-635, Urbana, IL 61801

Corresponding author: Eric Potash, epotash@illinois.edu

Abstract

Methods for measuring recreational water quality vary in analysis time, precision, availability, and cost. Decision-makers often use predictions from statistical models to compensate for the shortcomings of available measurements. However, model validation and comparison has largely omitted measurement error (defined here as variation in both sampling and the measurement technique) as an important source of uncertainty during validation. It is unknown how this omission affects estimates of model performance and comparisons between models. This study aims to fill this gap. First we derive the bias incurred when omitting measurement error in calculating a model's mean squared error. We then develop a non-parametric validation method to correct estimates of mean squared error. To study other metrics of prediction performance (mean absolute error, sensitivity, precision, etc.) we develop a second validation method that uses simulations from a Bayesian validation model. These methods are applied to a comparison of two prediction models (random forest and nearest neighbor) used to predict the level of fecal indicator bacteria at 9 recreational beaches in the city of Chicago. We find that accounting for measurement error significantly changes estimates of model performance. Moreover it reveals substantial uncertainty underlying some of these estimates.

1 Introduction

Recreational waterways are subject to contamination by bacteria from various sources including stormwater, sewage, and wildlife (Whitman & Nevers, 2008). To mitigate the public's exposure to contaminated water and associated gastrointestinal illness (Prüss, 1998), managers of recreational beaches monitor the presence of fecal indicator bacteria (FIB) as a proxy measure of contamination. Managers issue warnings or close sites based on this information. There is a trade off between the protective public health benefits of these actions and the recreational benefits of access to waterways (Rabinovici et al., 2004). A major challenge in this decision process is how to appropriately account for measurement error in FIB data, which can be substantial (Whitman & Nevers, 2004; Whitman et al., 2010). Here we define measurement error as the combined effect of error in the measurement process and in-situ sampling variability.

Measurement error has long been recognized as a major issue in water resources management, and the literature is rich with methods to incorporate measurement uncertainty in modeling and decision analysis. In hydrology, for example, Bayesian rainfall-runoff models have been developed to account for significant measurement error in catchment-scale precipitation to support improved parameter inference, predictive uncertainty bounds, and structural error diagnostics (Kuczera et al., 2006; Vrugt et al., 2008; Renard et al., 2011). Similar methods have also been extended to urban stormwater models to propagate bias and variance in both input (e.g. rainfall) and calibration (e.g., stormwater quality) data through the model fitting process (Dotto et al., 2014). Accommodations for measurement error have also been incorporated into decision-making processes, for instance with respect to groundwater remediation. For example, (Liu et al., 2012) used a value-of-information approach to estimate remediation cost reductions afforded by reduced model, parameter, and measurement uncertainty. Likewise, (Leube et al., 2012) used Bayesian methods to consider the effect of integrated groundwater modeling uncertainties (including measurement error) on optimal sampling design.

Measurement error has also played a prominent role in recreational water quality analysis. Modeling in this literature is often oriented towards decision support, where model-based predictions of FIB concentrations (including estimated moments or percentiles of measured data) are compared to water quality standards to guide

65 management actions. A significant body of work has considered the impacts of mea-
66 surement error on these decisions. For instance, several studies have used Bayesian
67 analyses to explore the potential of concentration-based FIB standards that account
68 for measurement error in indirect FIB concentration proxy measures (Gronewold et
69 al., 2008; Gronewold & Borsuk, 2010; Gronewold et al., 2017). A similar approach was
70 used to show that a significant fraction of space-time variability in FIB proxy measures
71 is driven by errors in measurement techniques and not underlying variability of in-situ
72 FIB concentrations (Gronewold et al., 2013).

73 When trying to improve water quality management decisions in the presence of
74 model structural uncertainty, it is also common to compare the predictive performance
75 of multiple FIB concentration models. In this facet of recreational water quality mod-
76 eling, however, measurement error has been given less attention. When validating
77 prediction models, we found that researchers often ignored measurement error, simply
78 assuming that a measurement (or mean of multiple measurements) represented the true
79 bacteria level (Nevers & Whitman, 2011; Francy, 2013; Shively et al., 2016; Lucius et
80 al., 2019). This is true even in studies that consider measurement error in the model
81 estimation process (e.g., see figure 5 and associated discussion in Gronewold et al.
82 (2011)). The omission of measurement error thus distorts a comparison of prediction
83 performance across models, although the magnitude of this effect is unknown.

84 Given the methodological gap above, this study contributes two methods for
85 model validation that account for measurement error when evaluating and comparing
86 the performance of prediction models. The first is a non-parametric method that makes
87 minimal assumptions but is limited to a single metric of model performance, namely
88 mean squared error (MSE). The second is a Bayesian method that uses simulation from
89 the posterior distribution of a Bayesian measurement error model. This method has
90 the advantage of being applicable to any metric of model performance, including those
91 assessing the utility of predictions for decision-making around management-relevant
92 FIB thresholds.

93 These methods are generally applicable to any inter-model comparison, and are
94 thus relevant across a range of modeling exercises in water quantity and quality anal-
95 ysis, not to mention other domains. However, they are particularly relevant to recre-
96 ational water quality modeling given the common task of comparing multiple FIB
97 concentration models for decision support and the high degree of measurement er-
98 ror in these data. We thus demonstrate the approach in a case study of recreational
99 beaches in Chicago, which has been used extensively to compare statistical models
100 that aid in prediction of bacteria levels (Nevers & Whitman, 2011; Shively et al., 2016;
101 Lucius et al., 2019).

102 2 Materials and methods

103 In our analysis we present both prediction models (section 2.2) and validation
104 methods (2.3). The prediction models are used to predict FIB levels at unsampled
105 beaches to support beach management decisions. The validation methods are used
106 to evaluate and compare these prediction models and their resulting decisions. The
107 focus of this study is on the methods for validation, not the specific prediction models
108 being validated. We compare a commonly used method for validation (termed naive
109 validation) against two new methods (a non-parametric method and Bayesian method)
110 that account for measurement error. We note that Bayesian validation relies on an
111 auxiliary model (termed the Bayesian validation model), which is used strictly to
112 validate other prediction models and not for prediction itself (more detail given in
113 section 2.3.4).

2.1 Study site and data

The city of Chicago has 23 beaches along approximately 42 km of the Southwest shoreline of Lake Michigan. Of these, 19 beaches (figure 1) are currently subject to FIB monitoring during the swimming season from late May to early September. The beaches receive about 20 million visits during this period each year (Nevers & Whitman, 2011).

Traditionally, administrators collected two samples per site for culture measurement of *E. coli* in terms of colony forming units (CFU) per 100 mL. Management decisions were made on the basis of the geometric mean measurement exceeding 235 CFU/100 mL. Culture measurements take at least 12-24 hours due to the bacteria incubation period. Because water quality can change rapidly, decisions based on measurements that are subject to such delays are likely to result in unnecessary closures as well as exposure (Kinzelman et al., 2003).

Starting in 2015, quantitative polymerase chain reaction (qPCR) measurements of *Enterococci* have been employed. This method quantifies indicator bacteria in less than two hours in terms of cell equivalents (CE) by comparing the sample to a calibrator with known number of *Enterococci* cells. A subset of these measurements are shown in figure 1. For details and comparisons of culture and qPCR measurements and their consequences see Noble et al. (2010), U.S. EPA (2012), and Dorevitch et al. (2017). Managers in Chicago currently estimate FIB levels at each beach each day using the geometric mean of two qPCR sample measurements. Management decisions are made, following U.S. EPA guidance (U.S. EPA, 2012), based on this estimate exceeding 1000 CE/mL.

Due to the cost of these measurements (Whitman et al., 2010) and the historical correlation of FIB levels between sites, the city has proposed reducing sampling to ten beaches and using a random forest (RF) model to predict levels at the remaining beaches (Lucius et al., 2019). The sampled beaches were chosen as follows. First, five beaches were selected to be sampled due to their historically high FIB levels. Next, the remaining beaches were grouped into five geographic clusters. Five beaches were selected to be sampled, one from each cluster. The five historically high FIB sites and five cluster representatives together give 10 sampled beaches, leaving 9 sites at which to make predictions (see figure 1). The prediction model uses daily meteorological and hydrological covariates collected between 2015-2019 for the months of May-September.

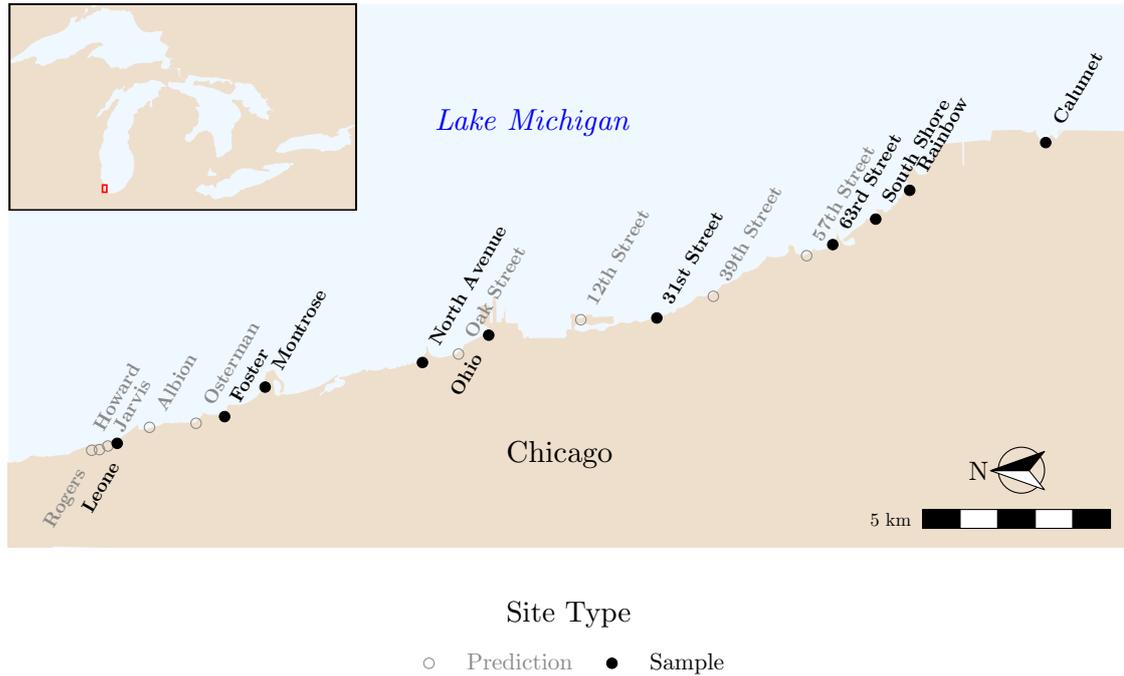
2.2 Prediction models of bacteria levels at unsampled sites

In this study we (re-)evaluate the predictions and management consequences of the RF model, and compare its performance against a benchmark NN model. The RF and NN models both serve as candidate *prediction* models. The purpose of this work is to assess how to validate and compare performance across prediction models given measurement error in the observations. We use the RF and NN prediction models to demonstrate our validation methods, but note that other prediction models could have been used for this purpose. Prior to describing the validation methods that are the focus of this work (see section 2.3), we first introduce notation and details of the specific predictive models used in the case study.

We denote the true (unobserved) level of *Enterococci* natural log cell equivalents per mL (log CE/mL) by θ_{jt} with $j = 1 \dots J$ a site index and $t = 1 \dots T$ a day index. Let Y_{ijt} be the observed measurements of θ_{jt} . In our case we typically have two measurements Y_{1jt}, Y_{2jt} which are replicates, taken at the same time and location.

Here we present NN-based and RF-based predictions of θ_{jt} at the prediction sites j . On day t , both predictions are based on the input vector of mean FIB measurements

(A) Map of study sites



(B) Measured fecal indicator bacteria levels

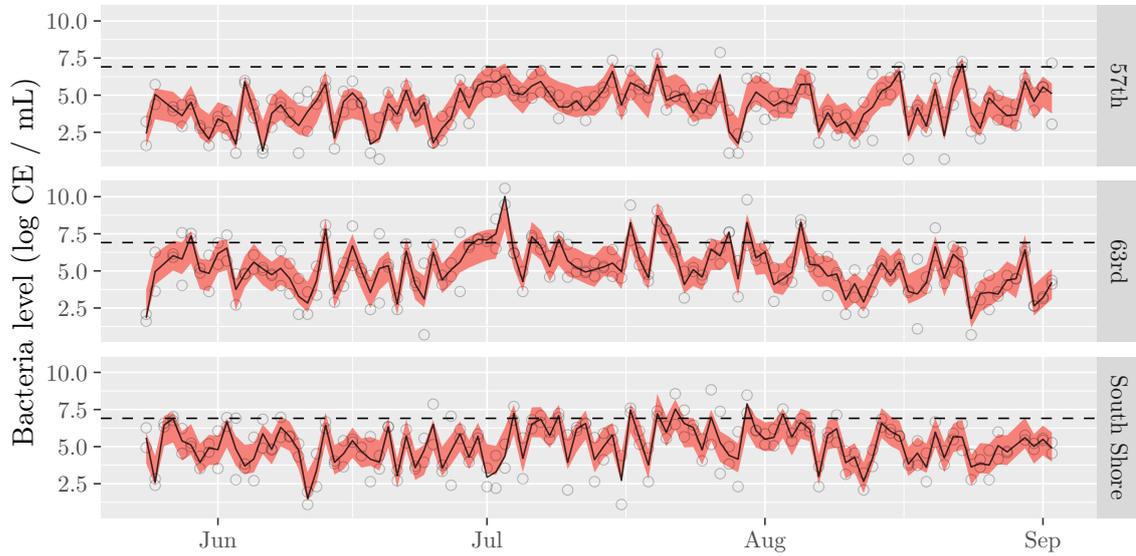


Figure 1: (A) Map of the 19 recreational beaches in Chicago on Lake Michigan (see inset for location within the Great Lakes) showing sample and prediction sites according to the proposed targeted sampling design of Lucius et al. (2019) described in section 2.2.2 and (B) daily fecal indicator bacteria levels during the 2019 beach season at three nearby sites. Gray circles are qPCR measurements (typically 2 per site per day), black line connects daily mean measurements, red region is 95% interval of Bayesian validation model posterior, and dashed line indicates action threshold of log 1000 cell equivalents (CE) per mL.

163 \bar{Y}_{jt} at the proposed ten sampled sites (figure 1). The RF prediction employs an
 164 additional input vector of K covariates X_{jt} varying by date and site. The outputs are
 165 predictions of the true FIB level at the proposed nine prediction sites.

166 **2.2.1 Nearest neighbor (NN) prediction model**

For prediction site j , let $n(j)$ be the index $(1, \dots, J)$ of the geographically nearest
 sampled site (figure 1A). The NN model predicts the FIB level at site j on date t to
 be equal to the mean level at this neighbor on the same date:

$$\hat{\theta}_{jt}^{\text{nn}} = \bar{Y}_{n(j),t}. \quad (1)$$

167 This model serves as a simple but practical benchmark.

168 **2.2.2 Random forest (RF) prediction model**

Lucius et al. (2019) proposed a “hybrid nowcast model” using a RF regression
 model with 400 trees (Breiman, 2001). The outcomes used to fit the model were the
 mean FIB levels at the prediction sites. The inputs to the model were the mean levels
 at the sampled sites together with covariates. Formally we can write the prediction
 model as a vector of functions

$$\hat{\theta}_{jt}^{\text{rf}}(\bar{Y}_t^{\text{sample}}, X_{jt}) \quad (2)$$

169 where X_{jt} is a vector of $K = 11$ covariates (varying by site and date) and $\bar{Y}_t^{\text{sample}}$ is
 170 the vector of average measurements at the ten sample sites on date t . For this study
 171 we refit the RF using our training set, which is larger than that of Lucius et al. (2019).
 172 The covariates X_{jt} mirror those of the original publication:

- 173 • Precipitation: 1 and 2-3 day total rainfall, 1-2 day change in water level
- 174 • Sunlight: 1 and 2-3 day average cloud cover
- 175 • Wind: 1 day average N/S/E/W wind speed
- 176 • Time: day of year, weekday indicator

177 where the weekday indicator is an indicator for whether the date is a weekday or
 178 weekend and 1 day, 1-2 day, and 2-3 day covariates aggregate over the period 24
 179 hours, 24-48 hours, and 48-72 hours prior.

180 **2.2.3 Calibration of exceedance predictions**

181 The above are prediction models of continuous FIB levels, but decisions are based
 182 on the binary event of exceeding 1000 CE/mL (section 2.1), for which an FIB level
 183 prediction must be transformed into a binary exceedance prediction. Some of our
 184 performance measures (e.g. precision) evaluate these exceedance predictions. For the
 185 baseline NN prediction model, we predicted an exceedance whenever the predicted FIB
 186 was greater than the 1000 CE/mL threshold. Since the RF is known to produce biased
 187 predictions, Lucius et al. (2019) calibrated a custom threshold so that the resulting
 188 specificity (equivalently false positive rate) matched that of a reference model. We
 189 follow this approach, taking NN as the reference model.

190 Exceedance could alternatively be predicted by modeling the binary outcome di-
 191 rectly, i.e. classification. However, we continue the standard practice in FIB prediction
 192 of modeling the continuous outcome, i.e. regression, as this uses all available informa-
 193 tion and allows us to use a single model for both continuous and binary outcomes.

194 **2.3 Validation methods for estimating prediction model performance**

195 The purpose of this study is to develop an approach to compare the performance
 196 of the above prediction models in the presence of measurement error. In validation
 197 we evaluate the fidelity of predicted states $\hat{\theta}$ to the true state θ by a function $L(\theta, \hat{\theta})$.
 198 Here L is one of various performance metrics (e.g. MSE) and θ and $\hat{\theta}$ are restricted
 199 to the prediction sites (figure 1) and dates t in a *test period* which we choose to be
 200 the most recent beach season, 2019. The RF model was fit using data from a *training*
 201 *period*, i.e. prior to 2019; the NN model does not require any fitting. We used a single
 202 training and test set, known as hold-out validation (Schneider & Moore, 2000), for
 203 simplicity. All of the methods below can be easily adapted for cross-validation with
 204 multiple folds.

205 Our challenge in validation is that we never observe θ_{jt} . In the literature, θ_{jt} is
 206 often assumed to be exactly equal to the mean measurement \bar{Y}_{jt} (Nevers & Whitman,
 207 2011; Francy, 2013; Shively et al., 2016; Lucius et al., 2019). Note that it is because
 208 these sites were in fact sampled that we can conduct this validation.

209 However, this method does not account for measurement uncertainty and it is
 210 unclear what the consequences of this omission are regarding the overall performance
 211 assessment of a prediction model or the comparison of multiple models. We term this
 212 method of validation *naive*, and propose two additional methods: *non-parametric* and
 213 *Bayesian*. The validation methods are described below and summarized in figure 2.

214 Note that there are two sources of variation accounting for the difference between
 215 the true FIB level θ_{jt} and an observation Y_{ijt} . The first is sampling variation due to
 216 the fact that a water sample is taken at a specific point in time and space (Whitman &
 217 Nevers, 2004). The second is measurement variation due to the qPCR technology used
 218 to analyze the sample (Whitman et al., 2010). As stated earlier, we define measurement
 219 error as the combination of these sampling and analysis variations.

220 **2.3.1 Naive validation not accounting for measurement error**

221 In naive validation we simply assume that the mean measurement is true: $\theta_{jt} =$
 222 \bar{Y}_{jt} . Then we can evaluate $L(\theta, \hat{\theta})$ as a point estimate of the performance (in contrast
 223 to non-parametric and Bayesian validation which estimate performance distributions).
 224 The use of naive validation is potentially flawed since the mean observation does not
 225 account for measurement error.

When the metric L is MSE, we can explicitly analyze the effect of measurement error. Assume a measurement error model

$$Y_{ijt} = \theta_{jt} + \epsilon_{ijt} \quad (3)$$

where ϵ_{ijt} are independent identically distributed measurement errors, defined as the sum of sampling and analysis errors. We do not assume a measurement error distribution but assume the errors have fixed (finite) variance $\text{Var}(\epsilon_{ijt}) = \tau^2$ independent of the measurement number i , site j , and date t . Then with \mathbb{E} denoting expectation over the random measurement errors ϵ we have

$$\mathbb{E}|\hat{\theta}_{jt} - \bar{Y}_{jt}|^2 = \mathbb{E}|\theta_{jt} + \bar{\epsilon}_{jt} - \hat{\theta}_{jt}|^2 \quad (4)$$

$$= \mathbb{E}[|\theta_{jt} - \hat{\theta}_{jt}|^2 + |\bar{\epsilon}_{jt}|^2 - 2\bar{\epsilon}_{jt}(\theta_{jt} - \hat{\theta}_{jt})] \quad (5)$$

$$= |\hat{\theta}_{jt} - \theta_{jt}|^2 + \frac{1}{2}\tau^2 \quad (6)$$

226 where we used the fact that the ϵ are independent of each other and both θ and $\hat{\theta}$ and
 227 there are two errors $\epsilon_{1jt}, \epsilon_{2jt}$ so that $\mathbb{E}[\bar{\epsilon}_{jt}^2] = \frac{1}{2}\tau^2$.

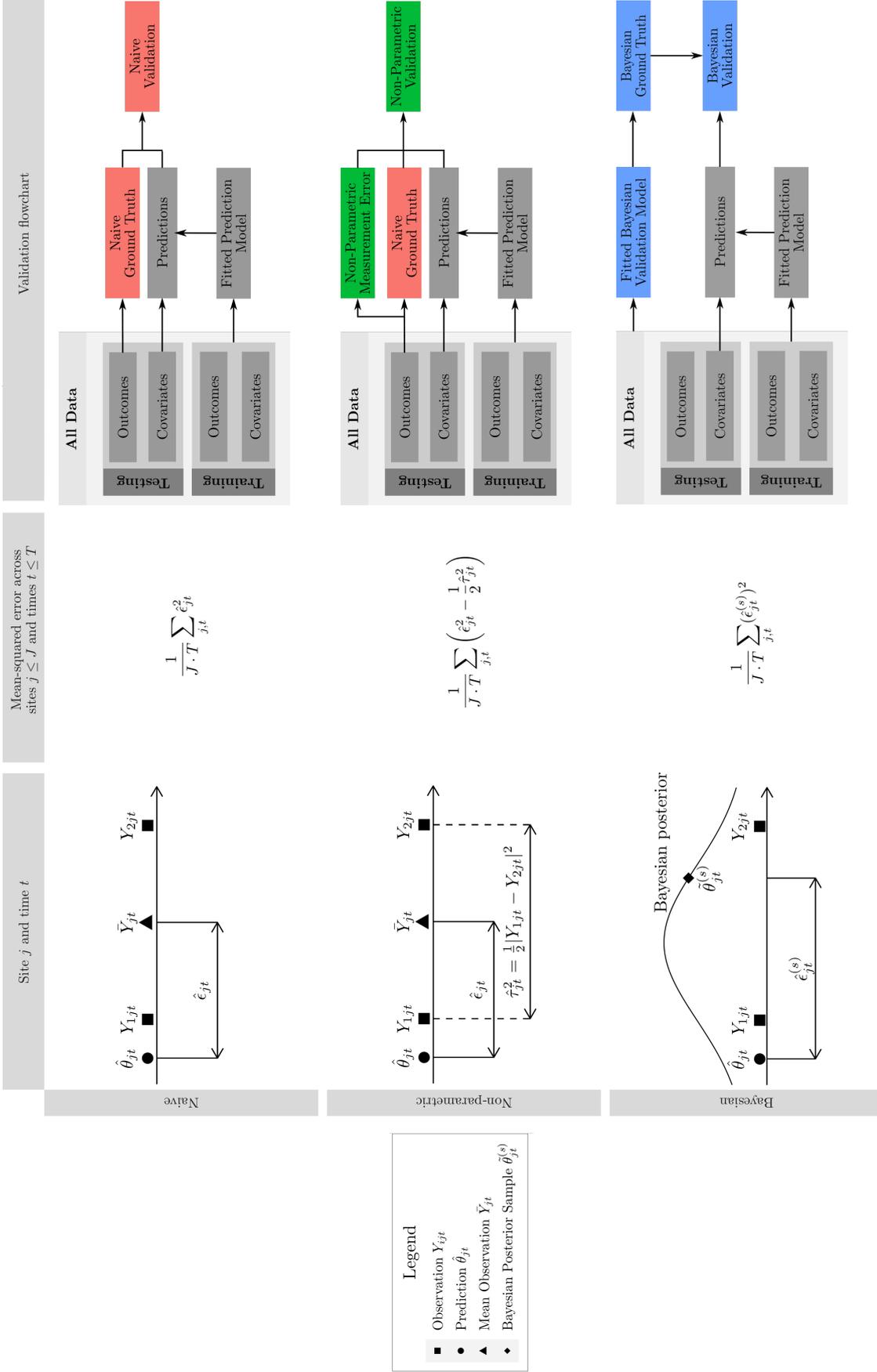


Figure 2: Validation methods. Prediction models are fit using training (2015-2018) covariates and outcomes, and use testing (2019) covariates to predict testing outcomes. For site j and time t the naive method simply compares the prediction $\hat{\theta}_{jt}$ with the mean observation \bar{Y}_{jt} to give a point estimate $\hat{\epsilon}_{jt} = \hat{\theta}_{jt} - \bar{Y}_{jt}$ of the error, which is averaged across sites j and times t to estimate MSE. The non-parametric method corrects the naive method by incorporating an estimate $\hat{\tau}$ of the measurement error. This is bootstrapped to estimate a sampling distribution. The Bayesian method compares the prediction to a sample of the posterior distribution of the FIB state in a Bayesian validation model that is fit to and conditioned on all the data. For each sample of the posterior $\hat{\theta}_{jt}^{(s)}$ (indexed by s) there is a corresponding sample of MSE. Abbreviations: mean squared error (MSE).

228 This formula means that under these mild assumptions, the naive estimate of
 229 MSE *overestimates* the true MSE by a multiple of the measurement error variance.
 230 Thus, the greater the measurement error variance τ^2 , the larger the distortion of naive
 231 validation for the particular metric of MSE.

However, since the distortion does not depend on which prediction model is being
 evaluated (e.g. RF or NN), the naive validation will give an unbiased estimate of the
 difference in performance, i.e.

$$\mathbb{E}[MSE(\bar{Y}, \hat{\theta}^{\text{rf}}) - MSE(\bar{Y}, \hat{\theta}^{\text{nn}})] = MSE(\theta, \hat{\theta}^{\text{rf}}) - MSE(\theta, \hat{\theta}^{\text{nn}}) \quad (7)$$

232 **2.3.2 Non-parametric validation accounting for measurement error**

233 Equation 4 shows that if we can estimate the measurement error τ then we can
 234 correct the bias of the naive estimate of MSE by subtracting it off. If we have more
 235 than one FIB measurement per beach-day in the data (as is the case in our data set),
 236 we can estimate τ using the standard sample variance estimator, and then average
 237 across beaches and days.

Namely if Y_{1jt} and Y_{2jt} are the two observations with measurement error ϵ_{1jt}
 and ϵ_{2jt} as in equation 3, we define the estimate

$$\hat{\tau}_{jt}^2 = \frac{1}{2} |Y_{1jt} - Y_{2jt}|^2 \quad (8)$$

which is unbiased because

$$\mathbb{E}[\hat{\tau}_{jt}^2] = \mathbb{E}\left[\frac{1}{2} |\epsilon_{1jt} - \epsilon_{2jt}|^2\right] \quad (9)$$

$$= \tau^2. \quad (10)$$

Combining 4 and 9 we have the following estimate for the mean-squared error of
 $\hat{\theta}_{jt}$:

$$\mathbb{E}[|\hat{\theta}_{jt} - \bar{Y}_{jt}|^2 - \frac{1}{2} \hat{\tau}_{jt}^2] = |\hat{\theta}_{jt} - \theta_{jt}|^2 \quad (11)$$

238 We average across sites j and dates t to estimate $MSE(\hat{\theta}, \theta)$. We bootstrap this
 239 estimate across t to estimate the sampling distribution (unlike naive validation which
 240 gives a point estimate).

241 We emphasize that this result does not make any distributional assumptions.
 242 We only assumed that the observations are equal to the true state plus independent
 243 measurement error (equation 3). However, this approach is limited to the specific error
 244 metric of MSE.

245 **2.3.3 Bayesian validation accounting for measurement error**

246 The approaches to validation presented above either: 1) assume the true FIB
 247 level θ_{jt} is equal to the mean of available observations (naive); or 2) indirectly estimate
 248 a specific error metric, MSE, under a specific sampling design (non-parametric). An
 249 alternative and more general approach is to use a Bayesian model to simulate the latent
 250 FIB state θ in validation. This is a form of Monte Carlo uncertainty propagation (ISO,
 251 2009).

252 We start with a Bayesian model (details of which are in section 2.3.4) of the FIB
 253 states and measurements given covariates $(\theta, Y|X)$ fit to all the data (i.e. covariates
 254 X and measurements Y). Then we sample θ at the prediction sites from the posterior
 255 distribution $\theta|Y, X$ conditioned on all the data. For convenience we write this posterior
 256 as $\tilde{\theta}$ and samples indexed by s as $\tilde{\theta}^{(s)}$. Then we can simply evaluate $L(\tilde{\theta}^{(s)}, \hat{\theta})$. Here, as

257 in naive and non-parametric validation, $\hat{\theta}$ are RF or NN predictions at the prediction
 258 sites based on covariates at those sites and measurements at the sampled sites. By
 259 repeatedly sampling $\tilde{\theta}^{(s)}$ and evaluating $L(\tilde{\theta}^{(s)}, \hat{\theta})$ we sample the target distribution
 260 $L(\theta, \hat{\theta})|Y, X$.

261 We emphasize that Y here includes all sites and times. That is, unlike the pre-
 262 diction models which only use observations from the sample sites during the testing
 263 period, the Bayesian simulations use measurements from the prediction sites them-
 264 selves (figure 1). In this way, the Bayesian model can support the validation of other
 265 prediction models, but should not be considered a prediction model itself.

266 Strictly speaking, the Bayesian validation method assumes that the Bayesian
 267 validation model is correct. However, it takes into account uncertainty in both model
 268 parameters and in observations due to measurement error. Moreover, the method is
 269 useful under the weaker assumption that the Bayesian simulations of θ provide a more
 270 realistic representation of the true FIB levels against which to compare the prediction
 271 models (and infer the target distribution $L(\theta, \hat{\theta})|Y, X$). Additionally, the MSE error
 272 metric inferred under Bayesian validation can be validated against the non-parametric
 273 approach, which makes fewer assumptions.

274 Yet compared to the non-parametric validation above, which only estimates MSE,
 275 Bayesian validation can be used to estimate any prediction performance metric. We
 276 consider several, including MSE, mean absolute error (MAE), and the area under the
 277 receiver operating curve (AUC) to evaluate predictions of the continuous FIB level.
 278 The remaining metrics (precision, sensitivity, specificity) use binary classifications that
 279 are obtained from continuous predictions using a threshold (see section 2.2.3). Preci-
 280 sion measures the proportion of exceedance predictions which are correct; sensitivity
 281 measures the proportion of exceedances which are correctly predicted; specificity mea-
 282 sures the proportion of non-exceedances which are correctly predicted. These latter
 283 metrics are particularly relevant to decision-making in recreational water quality man-
 284 agement, where binary decisions (e.g. site closure) are often based on water quality
 285 predictions exceeding a predetermined threshold.

286 **2.3.4 Bayesian validation model to facilitate Bayesian validation**

In order to implement Bayesian validation (section 2.3.3) we need a model of
 $\theta, Y|X$. Again, this model is not used for prediction, but rather uses all measurements
 (including those at prediction sites) to simulate the true FIB level θ and thus support
 validation of other prediction models that only use measurements at the sampling sites.
 We use a linear regression (on the log scale) model with coefficients varying by site:

$$\theta_{jt} = X_{jt}\beta_j + \eta_{jt} \quad (12)$$

287 where X_{jt} is a vector of K covariates (see section 2.2.2) and for each site j , β_j is
 288 a vector of K regression coefficients, and η_{jt} are model structural errors (distinct
 289 from measurement errors modeled below) capturing variation in the true FIB level
 290 not explained by the linear regression. We use the same covariates as in the RF (see
 291 section 2.2.2) but add an intercept and parameterize the day of year as a B-spline with
 292 4 degrees of freedom (since this model is linear as opposed to the non-linear RF). Thus
 293 $K = 15$.

On top of this regression we add three components. First we add a multivari-
 ate normal error distribution with covariance matrix Σ to model correlation in the
 structural errors across beaches on a given day t :

$$\eta_t \sim \text{Normal}(0, \Sigma) \quad (13)$$

294 This enables us to combine the measurements at other beaches with those at a given
 295 beach in estimating the bacteria level at that beach.

Second we add a multilevel structure on the coefficients, that is we have the second-level regression:

$$\beta_{jk} \sim \text{Normal}(Z\gamma_k, \sigma_{\beta_k}^2) \quad (14)$$

where Z is a $J \times L$ matrix of site-level covariates, γ_k is a vector of L second-level regression coefficients, and σ_{β_k} is the second-level residual standard deviation. Specifically we use $L = 3$ site level covariates including an intercept, breakwater length, and latitude (which is highly correlated with longitude, see figure 1). This second level regression allows the first level regression (equation 12) between FIB and the observation-level predictors (e.g. precipitation) to vary by site. Moreover, it allows us to partially pool information across sites and incorporate site-level information to more efficiently estimate the coefficients at a given beach (Stow et al., 2009; Cha et al., 2010). Incorporating site-level covariates also supports the (conditional) exchangeability of the sites in the model (Gelman et al., 2013).

The final component is an additive and normally distributed measurement error with variance τ^2 (Gronewold et al., 2009):

$$Y_{ijt} \sim \text{Normal}(\theta_{jt}, \tau^2). \quad (15)$$

Note that, unlike the RF model which is fit to beach-day mean levels \bar{Y}_{jt} , the Bayesian model is fit to the individual observations Y_{ijt} .

We put the following uninformative priors on these parameters (Gelman et al., 2013). Decomposing Σ into a correlation matrix Ω and a vector of coefficient scales σ

$$\Sigma = \text{diag}(\sigma) \cdot \Omega \cdot \text{diag}(\sigma) \quad (16)$$

we put a uniform prior over Ω and a Cauchy $_+(0, 1)$ prior on the components of σ . The second-level parameters γ_k and $\sigma_{\beta_k}^2$ are given uninformative Cauchy $(0, 1)$ and Cauchy $_+(0, 1)$ priors, respectively. All priors are defined after standardizing all predictors and the outcome.

We fit the Bayesian validation model using the Markov Chain Monte Carlo software Stan (Carpenter et al., 2017), which uses No-U-Turn sampling (Hoffman & Gelman, 2014), an extension of Hamiltonian Monte Carlo (Duane et al., 1987). We generated 4 chains with 1000 iterations each, saving the last 500 to produce $S = 2000$ samples from the joint posterior parameter distribution. We assessed mixing using the criteria $\hat{R} < 1.05$ and $n_{\text{eff}}/N > .001$ where \hat{R} is the Gelman-Rubin convergence statistic and n_{eff} is the effective sample size (Gelman et al., 2013).

With the uninformative prior on Σ we are making relatively weak assumptions about the covariance structure. This is possible in our application because of the relatively small number (19) of sites and the efficiency of Hamiltonian Monte Carlo. In applications with more sites, it may be useful to model the covariance in terms of the distance between sites j and k using a Gaussian process model or in terms of an adjacency matrix using a conditional autoregressive model (Gelfand et al., 2010).

3 Results

The Bayesian validation model was fit using 13,109 observations at the 19 beaches on 430 days between 2015 and 2019. The posterior distribution of measurement error variance τ^2 had median 0.77 (95% CI, 0.74 to 0.8). This was 30% of the variance of daily means \bar{Y}_{jt} of 2.5. Complete results of the fit are included in the supplementary material.

During the 2019 season 3,780 qPCR measurements were made over 102 days. The Bayesian posterior FIB level θ is displayed in figure 1. We restricted the test

333 period to those days with two samples at each of the 19 beaches. There were 67 such
 334 days.

335 Using the geometric mean FIB level on each beach-day, the median FIB level was
 336 92 CE/mL and 4.9% of these beach-days were in exceedance of the 1000 CE threshold.
 337 The Bayesian estimate of the median level was 93 CE (95% CI, 88 to 99 CE) with
 338 3.9% (95% CI, 3.1% to 4.7%) of beach days exceeded the threshold.

339 We first compare all three validation methods' estimates of MSE since this is
 340 the only metric where the non-parametric method is applicable. Then we compare
 341 Bayesian and naive estimates of all prediction performance metrics.

342 **3.1 Mean squared error under all validation methods**

343 We start by examining the three validation methods on the metric where they
 344 can all be compared, namely MSE. Figure 3(A) presents these estimates. While naive
 345 validation gives point estimates, non-parametric and Bayesian validation give distri-
 346 butions. Moreover, the latter give joint distributions of MSE for the two prediction
 347 models and so yield distributions for the difference in MSE between the two prediction
 348 models.

349 There are four findings here. First, we anticipated that naive validation would
 350 give a positively biased estimate of mean-squared error (4) and we see that it does
 351 give larger estimates than both non-parametric and Bayesian validation. For both RF
 352 and NN prediction models, the naive estimates of MSE lie above the 95% intervals
 353 estimated by non-parametric validation. Because non-parametric validation accounts
 354 for measurement error, we are inclined to trust its results and dismiss naive validation
 355 which is a priori flawed.

356 Second, we find as expected (equation 7) that while naive validation overstates
 357 the MSE of both prediction models, the estimated *difference* between the prediction
 358 models agrees with the difference given by non-parametric validation.

359 Third, we find remarkable agreement between non-parametric and Bayesian MSE
 360 estimates. Because non-parametric validation makes few assumptions, this agreement
 361 provides evidence to support our use of the Bayesian validation method to further
 362 explore the performance of prediction models and metrics for which we do not have a
 363 non-parametric method (discussed next).

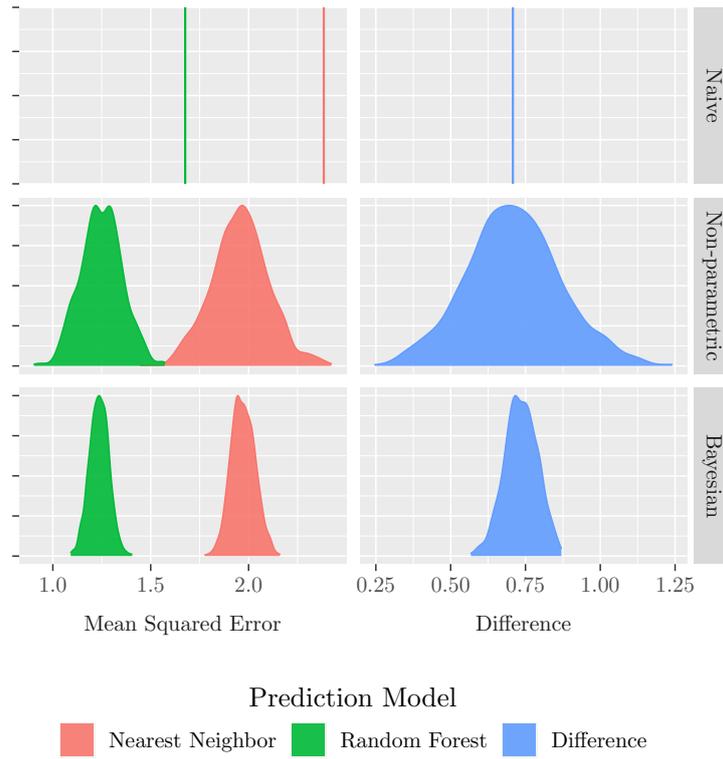
364 Fourth, the Bayesian method provides narrower uncertainty around its estimates
 365 than the non-parametric method. This may be explained by the fact that when esti-
 366 mating the squared error of a given prediction $\hat{\theta}_{jt}$, the non-parametric method only
 367 uses the measurements Y_{ijt} at the site while the Bayesian method uses the additional
 368 information of covariates and measurements at other sites.

369 **3.2 All performance measures under naive and Bayesian validation**

370 We proceed to evaluate the full set of performance metrics using Bayesian and
 371 naive validation methods. We started by using Bayesian validation to estimate the
 372 expected sensitivity of the NN prediction model with a binary classification threshold
 373 of 1000 CE. The estimate was 95.6%, and to match this (section 2.2.3), a threshold of
 374 440 for the RF was calibrated. Estimates for all prediction performance metrics are
 375 shown in figure 3(B).

376 According to both naive and Bayesian validation, RF outperforms NN in all
 377 metrics. However, the discrepancy between Bayesian and naive validation, first docu-
 378 mented for MSE in section 3.1 above, continues here across more metrics. Unlike
 379 MSE which was systematically overestimated (i.e. pessimistic) using naive validation,

(A) Mean squared error in all validation methods



(B) All performance metrics in naive and Bayesian validation

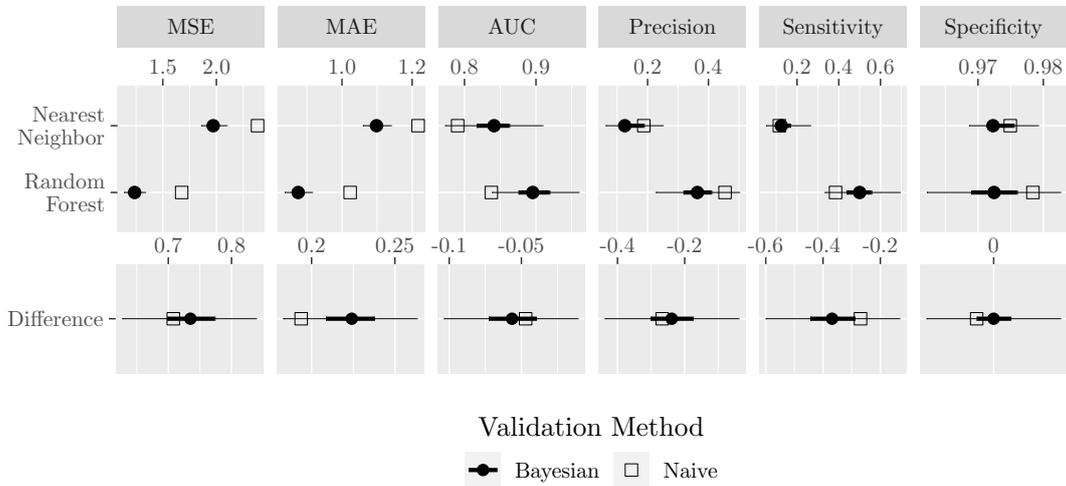


Figure 3: (A) Mean squared error of nearest neighbor and random forest predictions and their difference estimated using naive, non-parametric, and Bayesian validation methods. (B) Prediction performance estimates using naive and Bayesian validation. Solid dots and intervals show median and 50% and 95% credible intervals, respectively, using simulation to account for measurement error. Open squares show naive estimates without accounting for measurement error. Abbreviations: mean squared error (MSE), mean absolute error (MAE), area under receiver operating characteristic curve (AUC).

380 other measures are variously pessimistic and optimistic, such as when naive validation
 381 estimates the precision of the RF to be 0.45 while the Bayesian point estimate is 0.36.

382 However, as for MSE above, for each metric the bias of naive validation (relative
 383 to Bayesian validation) is the same for both RF and NN. Naive estimates of the
 384 *difference* in prediction performance benefit from this consistency in that the biases of the
 385 difference estimates are not larger than those of the absolute estimates. However,
 386 the discrepancy between naive and Bayesian validation of these estimated differences
 387 can still be quite large, as for example sensitivity where naive validation estimates RF
 388 to be an improvement of 0.27 while the Bayesian estimate is 0.37.

389 The discrepancy between naive and Bayesian validation turns out to be greatest
 390 for MSE and MAE which measure performance of continuous FIB level predictions,
 391 and less for the others which measure binary prediction performance.

392 While the Bayesian estimates of MSE and MAE are relatively precise, the es-
 393 timates of precision and sensitivity are very uncertain. This is explained by the fact
 394 that exceedances (positive events) form the denominator in the sensitivity metric and
 395 contribute to the numerator in precision. As a tail event defined by a sharp threshold,
 396 exceedance is difficult to measure and model precisely.

397 These metrics as well as their uncertainty capture the efficacy of the predictions
 398 for binary management decision making such as issuing a swim advisory or closing a
 399 site. Sensitivity for example is the proportion of elevated FIB events that are correctly
 400 predicted (and hence acted on), which for the RF is 0.50 (95% CI, 0.33 to 0.69).
 401 Precision on the other hand is the proportion of predicted elevated FIB events (hence
 402 actions taken) that are correct, which for the RF is 0.36 (95% CI, 0.23 to 0.50).

403 4 Conclusion

404 The omission of measurement error, specifically in validation, is ubiquitous in
 405 the water resources literature (e.g. Dawson and Wilby (2001); Berenguer et al. (2005);
 406 Biondi et al. (2012); Lohani et al. (2012); Shortridge et al. (2016)), including prediction
 407 models for recreational water quality (e.g. Nevers and Whitman (2011); Francy (2013);
 408 Shively et al. (2016); Lucius et al. (2019)). In this technical note we examined the
 409 effect of this omission.

410 For the specific prediction performance metric of MSE we showed that ignoring
 411 measurement error biases validation results (equation 4). The size of the bias depends
 412 on the size of the measurement error, which is very large in our context of recre-
 413 ational water quality. Next we contributed two new methods for model validation and
 414 inter-comparison that account for measurement error. The first was a non-parametric
 415 method making few assumptions but limited to the metric of MSE. The second was a
 416 Bayesian method that uses simulations from a parametric model to estimate any per-
 417 formance metric. We applied these methods to the evaluation of prediction models of
 418 FIB levels at beaches in Chicago and found that not accounting for measurement error
 419 significantly mis-estimated model performance across a range of metrics. Moreover it
 420 failed to quantify the uncertainty of prediction performance. Our non-parametric and
 421 Bayesian approaches overcame these issues.

422 Accurate model skill assessments are important. These estimates are required
 423 by water quality managers to understand the utility of model predictions for decision-
 424 making. Bias in estimated performance metrics could skew how decision-makers inter-
 425 pret model predictions or select among competing models, as could the presentation
 426 (or lack thereof) of performance uncertainty. More generally, performance estimates
 427 and their uncertainty are essential to understanding the public health consequences of
 428 management decisions made on the basis of these models. Measures of model perfor-

mance are also used to inform decisions about additional sampling (if deemed necessary to improve performance), which could be costly. For example, if the city of Chicago were to conclude based on an assessment of model performance that they needed two samples per beach-day at the 19 beaches, rather than the current proposal which includes roughly half the number of samples, the additional sampling cost would total about \$57,000 per season assuming an estimated qPCR analysis cost of \$30 per sample (Griffith & Weisberg, 2011).

Both the non-parametric and Bayesian approaches to validation proposed in this study help overcome limitations of a naive approach. However, the Bayesian approach to validation is more flexible in its ability to compare models across a variety of metrics, some of which might be particularly relevant to decision-making (e.g., predictions of exceedances over a water quality threshold).

Additional developments could improve the Bayesian model. Temporal autocorrelation was considered but initial testing (not shown) confirmed previous findings that system dynamics are too fast (Dorevitch et al., 2017). The regression coefficients β_{jk} could be modeled as correlated using a multivariate normal distribution. A more sophisticated approach to spatial correlation is also possible, e.g. using an adjacency matrix within a conditional autoregressive model (Gelfand et al., 2010). There is some evidence suggesting that measurement error may vary with the bacteria level (Whitman et al., 2010), which could be modeled using a heteroskedastic measurement error. While the measurements in our dataset are only at a single location and time for each beach and date, it has been shown that there is substantial variation spatially within each beach and temporally within each day (Whitman & Nevers, 2004; Boehm, 2007). With the relevant data, our model could be extended to these finer scales. In applications where predictions models produce probabilistic forecasts, the Bayesian validation method could be further developed with performance measures that compare the forecast distribution with the Bayesian posterior. These efforts are left for future work. Importantly though, even if the Bayesian model is not the best prediction model of FIB levels (as compared to, say, a machine learning model), it enables us to incorporate all available information into simulations of uncertain bacteria concentrations. This study shows how those simulations can be used to validate prediction models for more realistic assessments of skill compared to a naive approach.

Acknowledgments

This research did not receive any specific grant from funding agencies in the public, commercial, or not-for-profit sectors. Thanks to Dan Black, Steven Durlauf, Jeff Johnston, Andrew Gelman, Jim Savage, Jackie Shadlen, and Rob Trangucci for useful conversations. Thanks also to Lucius et al. for transparency about their model and data and assistance in replicating their results. And thanks to the associate editor and reviewers for helpful feedback.

Datasets used in this research are available from the following websites: bacteria test results (<https://data.cityofchicago.org/>); site-specific rainfall, cloud cover, and wind speed data (<https://darksky.net>); and Lake Michigan water levels at Calumet Harbor (<https://tidesandcurrents.noaa.gov/>).

References

- Berenguer, M., Corral, C., Sánchez-Diezma, R., & Sempere-Torres, D. (2005). Hydrological validation of a radar-based nowcasting technique. *Journal of Hydrometeorology*, 6(4), 532–549.
- Biondi, D., Freni, G., Iacobellis, V., Mascaro, G., & Montanari, A. (2012). Validation of hydrological models: Conceptual basis, methodological approaches and

- 478 a proposal for a code of practice. *Physics and Chemistry of the Earth, Parts*
 479 *A/B/C*, 42, 70–76.
- 480 Boehm, A. (2007). Enterococci concentrations in diverse coastal environments ex-
 481 hibit extreme variability. *Environmental Science & Technology*, 41(24), 8227–
 482 8232.
- 483 Breiman, L. (2001). Random forests. *Machine learning*, 45(1), 5–32.
- 484 Carpenter, B., Gelman, A., Hoffman, M. D., Lee, D., Goodrich, B., Betancourt, M.,
 485 ... Riddell, A. (2017). Stan: A probabilistic programming language. *Journal*
 486 *of statistical software*, 76(1).
- 487 Cha, Y., Stow, C. A., Reckhow, K. H., DeMarchi, C., & Johengen, T. H. (2010).
 488 Phosphorus load estimation in the Saginaw River, MI using a Bayesian hierar-
 489 chical/multilevel model. *Water Research*, 44(10), 3270–3282.
- 490 Dawson, C., & Wilby, R. (2001). Hydrological modelling using artificial neural net-
 491 works. *Progress in Physical Geography*, 25(1), 80–108.
- 492 Dorevitch, S., Shrestha, A., DeFlorio-Barker, S., Breitenbach, C., & Heimler, I.
 493 (2017). Monitoring urban beaches with qPCR vs. culture measures of fecal
 494 indicator bacteria: Implications for public notification. *Environmental Health*,
 495 16(1), 45.
- 496 Dotto, C. B. S., Kleidorfer, M., Deletic, A., Rauch, W., & McCarthy, D. T. (2014).
 497 Impacts of measured data uncertainty on urban stormwater models. *Journal of*
 498 *Hydrology*, 508, 28–42.
- 499 Duane, S., Kennedy, A. D., Pendleton, B. J., & Roweth, D. (1987). Hybrid Monte
 500 Carlo. *Physics letters B*, 195(2), 216–222.
- 501 Franczy, D. S. (2013). *Developing and implementing predictive models for estimating*
 502 *recreational water quality at Great Lakes beaches*. US Department of the Inter-
 503 ior, US Geological Survey.
- 504 Gelfand, A. E., Diggle, P., Guttorp, P., & Fuentes, M. (2010). *Handbook of spatial*
 505 *statistics*. CRC press.
- 506 Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B.
 507 (2013). *Bayesian data analysis*. Chapman and Hall/CRC.
- 508 Griffith, J. F., & Weisberg, S. B. (2011). Challenges in implementing new technology
 509 for beach water quality monitoring: lessons from a california demonstration
 510 project. *Marine Technology Society Journal*, 45(2), 65–73.
- 511 Gronewold, A. D., Borsuk, M., Wolpert, R., & Reckhow, K. (2008). An assessment
 512 of fecal indicator bacteria-based water quality standards. *Environmental Sci-*
 513 *ence & Technology*, 42(13), 4676–4682.
- 514 Gronewold, A. D., & Borsuk, M. E. (2010). Improving water quality assessments
 515 through a hierarchical Bayesian analysis of variability. *Environmental Science*
 516 *& Technology*, 44(20), 7858–7864.
- 517 Gronewold, A. D., Myers, L., Swall, J. L., & Noble, R. T. (2011). Addressing un-
 518 certainty in fecal indicator bacteria dark inactivation rates. *Water Research*,
 519 45(2), 652–664.
- 520 Gronewold, A. D., Qian, S. S., Wolpert, R. L., & Reckhow, K. H. (2009). Calibrat-
 521 ing and validating bacterial water quality models: A Bayesian approach. *Wa-*
 522 *ter Research*, 43(10), 2688–2698.
- 523 Gronewold, A. D., Sobsey, M. D., & McMahan, L. (2017). The compartment bag
 524 test (CBT) for enumerating fecal indicator bacteria: basis for design and inter-
 525 pretation of results. *Science of the Total Environment*, 587, 102–107.
- 526 Gronewold, A. D., Stow, C. A., Vijayavel, K., Moynihan, M. A., & Kashian, D. R.
 527 (2013). Differentiating enterococcus concentration spatial, temporal, and ana-
 528 lytical variability in recreational waters. *Water Research*, 47(7), 2141–2152.
- 529 Hoffman, M. D., & Gelman, A. (2014). The No-U-Turn sampler: adaptively set-
 530 ting path lengths in Hamiltonian Monte Carlo. *Journal of Machine Learning*
 531 *Research*, 15(1), 1593–1623.

- 532 ISO. (2009). *Uncertainty of measurement — part 1: Introduction to the expression*
 533 *of uncertainty in measurement* (Vol. 98-1; Tech. Rep.). Geneva, CH.
- 534 Kinzelman, J., Ng, C., Jackson, E., Gradus, S., & Bagley, R. (2003). Enterococci
 535 as indicators of Lake Michigan recreational water quality: comparison of two
 536 methodologies and their impacts on public health regulatory events. *Appl.*
 537 *Environ. Microbiol.*, *69*(1), 92–96.
- 538 Kuczera, G., Kavetski, D., Franks, S., & Thyer, M. (2006). Towards a Bayesian total
 539 error analysis of conceptual rainfall-runoff models: Characterising model error
 540 using storm-dependent parameters. *Journal of Hydrology*, *331*(1-2), 161–177.
- 541 Leube, P., Geiges, A., & Nowak, W. (2012). Bayesian assessment of the expected
 542 data impact on prediction confidence in optimal sampling design. *Water Re-*
 543 *sources Research*, *48*(2).
- 544 Liu, X., Lee, J., Kitanidis, P. K., Parker, J., & Kim, U. (2012). Value of information
 545 as a context-specific measure of uncertainty in groundwater remediation. *Wa-*
 546 *ter Resources Management*, *26*(6), 1513–1535.
- 547 Lohani, A., Kumar, R., & Singh, R. (2012). Hydrological time series modeling: A
 548 comparison between adaptive neuro-fuzzy, neural network and autoregressive
 549 techniques. *Journal of Hydrology*, *442*, 23–35.
- 550 Lucius, N., Rose, K., Osborn, C., Sweeney, M. E., Chesak, R., Beslow, S., &
 551 Schenk Jr, T. (2019). Predicting *E. coli* concentrations using limited qPCR
 552 deployments at Chicago beaches. *Water Research X*, *2*, 100016.
- 553 Nevers, M. B., & Whitman, R. L. (2011). Efficacy of monitoring and empirical
 554 predictive modeling at improving public health protection at Chicago beaches.
 555 *Water Research*, *45*(4), 1659–1668.
- 556 Noble, R. T., Blackwood, A. D., Griffith, J. F., McGee, C. D., & Weisberg, S. B.
 557 (2010). Comparison of rapid quantitative PCR-based and conventional culture-
 558 based methods for enumeration of *Enterococcus* spp. and *Escherichia coli* in
 559 recreational waters. *Appl. Environ. Microbiol.*, *76*(22), 7437–7443.
- 560 Prüss, A. (1998). Review of epidemiological studies on health effects from exposure
 561 to recreational water. *International Journal of Epidemiology*, *27*(1), 1–9.
- 562 Rabinovici, S. J., Bernknopf, R. L., Wein, A. M., Coursey, D. L., & Whitman, R. L.
 563 (2004). Economic and health risk trade-offs of swim closures at a lake michigan
 564 beach. *Environmental Science & Technology*, *38*(10), 2737–2745.
- 565 Renard, B., Kavetski, D., Leblois, E., Thyer, M., Kuczera, G., & Franks, S. W.
 566 (2011). Toward a reliable decomposition of predictive uncertainty in hydro-
 567 logical modeling: Characterizing rainfall errors using conditional simulation.
 568 *Water Resources Research*, *47*(11).
- 569 Schneider, J., & Moore, A. (2000, February). *A locally weighted learning tutorial us-*
 570 *ing Vizier 1.0* (Tech. Rep. No. CMU-RI-TR-00-18). Pittsburgh, PA: Carnegie
 571 Mellon University.
- 572 Shively, D. A., Nevers, M. B., Breitenbach, C., Phanikumar, M. S., Przybyla-Kelly,
 573 K., Spoljaric, A. M., & Whitman, R. L. (2016). Prototypic automated contin-
 574 uous recreational water quality monitoring of nine Chicago beaches. *Journal of*
 575 *Environmental Management*, *166*, 285–293.
- 576 Shortridge, J. E., Guikema, S. D., & Zaitchik, B. F. (2016). Machine learning
 577 methods for empirical streamflow simulation: a comparison of model accuracy,
 578 interpretability, and uncertainty in seasonal watersheds. *Hydrology & Earth*
 579 *System Sciences*, *20*(7).
- 580 Stow, C. A., Lamon, E. C., Qian, S. S., Soranno, P. A., & Reckhow, K. H. (2009).
 581 Bayesian hierarchical/multilevel models for inference and prediction using
 582 cross-system lake data. In *Real World Ecology* (pp. 111–136). Springer.
- 583 U.S. EPA. (2012). *Recreational water quality criteria* (Tech. Rep.).
- 584 Vrugt, J. A., Ter Braak, C. J., Clark, M. P., Hyman, J. M., & Robinson, B. A.
 585 (2008). Treatment of input uncertainty in hydrologic modeling: Doing hydrolog-
 586 ical backward with Markov chain Monte Carlo simulation. *Water Resources*

- 587 *Research*, 44(12).
588 Whitman, R. L., Ge, Z., Nevers, M. B., Boehm, A. B., Chern, E. C., Haugland,
589 R. A., . . . others (2010). Relationship and variation of qPCR and culturable
590 Enterococci estimates in ambient surface waters are predictable. *Environmen-*
591 *tal Science & Technology*, 44(13), 5049–5054.
592 Whitman, R. L., & Nevers, M. B. (2004). *Escherichia coli sampling reliability at*
593 *a frequently closed Chicago beach: monitoring and management implications*.
594 ACS Publications.
595 Whitman, R. L., & Nevers, M. B. (2008). Summer E. coli patterns and responses
596 along 23 Chicago beaches. *Environmental Science & Technology*, 42(24), 9217–
597 9224.