

## Estimating changes in river faecal coliform loading using nonparametric multiplicative regression

Christopher J. Schulz and Gary W. Childers

### ABSTRACT

Faecal coliform (FC) concentration was monitored weekly in the Tangipahoa River over an eight year period. Available USGS discharge and precipitation data were used to construct a nonparametric multiplicative regression (NPMR) model for both forecasting and backcasting of FC density. NPMR backcasting and forecasting of FC allowed for estimation of concentration for any flow regime. During this study a remediation effort was undertaken to improve disinfection systems of contributing municipal waste water treatment plants in the watershed. Time-series analysis of FC concentrations demonstrated a drop in FC levels coinciding with remediation efforts. The NPMR model suggested the reduction in FC levels was not due to climate variance (i.e. discharge and precipitation changes) alone. Use of the NPMR method circumvented the need for construction of a more complex physical watershed model to estimate FC loading in the river. This method can be used to detect and estimate new discharge impacts, or forecast daily FC estimates.

**Key words** | indicators, pathogens, pollution modelling, public health, risk assessment, water criteria

**Christopher J. Schulz**  
**Gary W. Childers** (corresponding author)  
Department of Biological Sciences,  
Southeastern Louisiana University,  
Hammond, LA 70402,  
USA  
Tel.: +1 (985)549-3503  
Fax: +1 (985)549-3851  
E-mail: [gchilders@selu.edu](mailto:gchilders@selu.edu)

### INTRODUCTION

The importance of surface water monitoring cannot be understated as many diseases are acquired from contact with contaminated water bodies (Hurst & Murphy 1996). The presence of faecal pollution indicators in natural waters provides an index by which the public health risk can be gauged, in part because of observed relations to pathogens (Arvanitidou *et al.* 1997; Payment *et al.* 2000), and elevated health risk (Prüss 1998; Wade *et al.* 2003). Regulatory bodies have adopted bacteriological water standards under the guidance of the EPA (Dufour 1984) based on incidence of gastrointestinal illness. Theoretically, if the bacteriological water quality standards such as faecal indicator bacteria (FIB) are monitored continuously and in real time, then public health risks can be lowered by reduction of exposure to pathogens. However, there are major deficiencies with such approaches, including: FIB are expensive to monitor in

high frequency; there is a time lag between sample acquisition and reporting due to methodological limitations; and environmental conditions can change from the time of sampling (Leclerc *et al.* 2001; Karine Lemarchand 2003; Horman *et al.* 2004; Harwood *et al.* 2005). In general, FIB monitoring of water bodies on a weekly or monthly basis is an important tool for identifying impaired waters, but fails to provide information that is relevant on daily timescales. It has been suggested that statistical models can fill the gap in data and provide 'now-casting' for faecal pollution indicators (Boehm *et al.* 2009). The inspiration of this work was to investigate a novel statistical approach allowing for incorporation of readily available stream-flow and precipitation data to reduce uncertainty in FIB modelling.

An alternative to daily monitoring of FIB is forecast modelling. Optimally forecast models can utilize environmental

predictors that are readily assessable via continuous monitoring, such as USGS real-time hydrological data. The techniques used to derive a predictive FIB model are varied, but can be generalized as mechanistic models (physical), data driven statistical approaches (empirical), or mixtures of the two (Mahloch 1974). Physical models require formulation of relationships between environmental predictors and FIB loading *a priori*. General physical modelling structures are freely available from the EPA, such as Water Quality Analysis Simulation Program (WASP6.1) and River and Stream Quality Model (Qual2K). Physical modelling requires that the appropriate input variables are known and their interactions must be explicitly formulated, albeit often in an iterative fashion, to produce an acceptable predictive tool. This approach can be accomplished by addition of a pathogen transport model to a hydrologic model. A physical model (Kashefipour *et al.* 2002; Steets & Holden 2003) can improve fit with data through incorporation of additional assumptions (such as bacterial decay constants). The resources needed to derive an explicit determinative model can become limiting, especially considering the expertise needed to produce a reliable hydrological model.

Although physically based modelling techniques are available, the implementation has yet to become common practice (Jamieson *et al.* 2004). A likely reason for such is that each drainage basin has unique characteristics, such as, but not limited to: climate variation, soil properties, drainage area, roughness of river bed, hydraulic radius, land use, tidal influences and discharge. Proper incorporation of the environmental information and parameter estimation is therefore not a routine task. Alternatively, purely empirical approaches seek to derive relationships given the data and make no *a priori* assumptions of the mechanisms responsible for the observed responses. An empirical approach can be a pragmatic modelling solution circumventing some of the shortcomings of physical modelling, especially when limited resources are available.

Previous empirical FIB models have utilized ordinary least squares (OLS) regression (Mahloch 1974; Eleria & Vogel 2005; McCarthy *et al.* 2007), logistic regression (Whitman *et al.* 2004; Eleria & Vogel 2005; Chandramouli *et al.* 2008) and artificial neural networks (ANN) (Brion & Lingireddy 1999; Chandramouli *et al.* 2007, 2008; Mas & Ahlfeld 2007; He & He 2008). Of these methods, OLS

regression and logistic regression methods are the simplest to implement and therefore are attractive tools for recreational water management. OLS regression and logistic regression attempt to fit a parametric response of FIB to predictors, either linear or sigmoidal, respectively. In the case of ANN, the response form is unclear, but a good overall model fit can be achieved (Lin *et al.* 2003; Mas & Ahlfeld 2007; He & He 2008). More recently, generalized additive modelling (GAM) of environmental factors associated with *Salmonella* occurrence has been implemented (Setti *et al.* 2009).

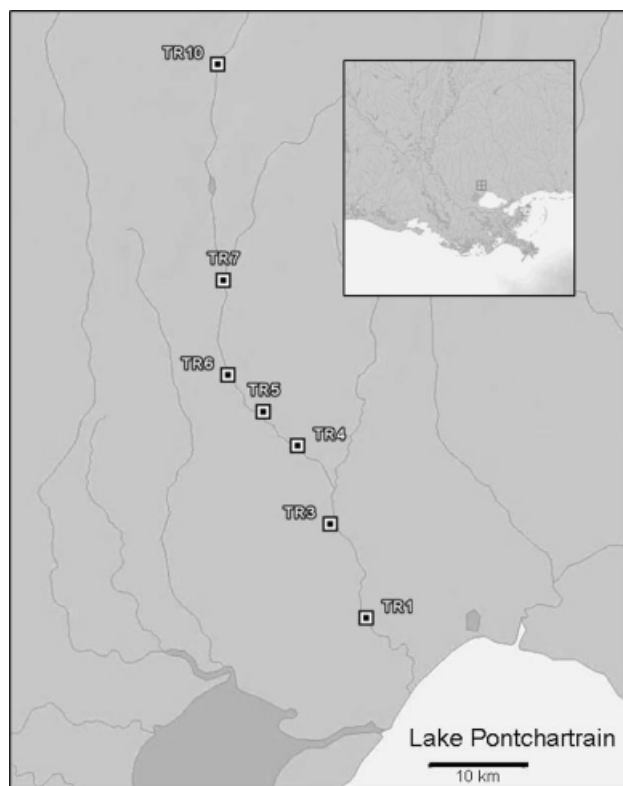
In this study, we have employed a new approach using nonparametric regression with interactions accounted for by multiplicative combining of predictors. The nonparametric multiplicative regression (NPMR) approach is a non-linear, nonparametric technique. A drawback of NPMR modelling is that no explicit formula is generated, requiring that data used to generate the model is also utilized in forecasting (or backcasting). This approach can readily accommodate response surface construction, allowing for visualization of predictor interactions. The NPMR method also has the benefit of ease of use, not requiring sophisticated model development through trial and error or mechanistic modelling techniques. This approach has recently been adapted for use in ecological habitat modelling (McCune 2006a, b), but to our knowledge, has not previously been applied to water quality modelling.

The NPMR approach shares roots with generalized additive models (GAM). A GAM differs from NPMR in that individual terms are summed with each term smoothed by a nonparametric function, whereas NPMR combines all predictors and weights multiplicatively. It is possible to construct a GAM model of the same structure as a NPMR model by limiting the number of terms to one and combining all predictors multiplicatively; thus, NPMR is a special case of a GAM (McCune 2006b). Even though it is possible to construct a NPMR model using the GAM framework, it would seem inappropriate to call a NPMR model a GAM model, as there would be no additive terms. This study investigated the use of NPMR for water quality modelling with respect to faecal coliform (FC) forecasting and backcasting to monitor changes in water quality. The goal of this study was to utilize NPMR as a method for evaluating water quality as a forecasting tool.

## MATERIALS AND METHODS

### Study site description

Surface water samples were collected from the Tangipahoa River in Louisiana (Figure 1) at site TR3. The river extends north to south for nearly 106 km (66 miles) through south-east Louisiana, drains 1,673 km<sup>2</sup> (646 square miles), and empties into Lake Pontchartrain. The Tangipahoa River is heavily utilized during summer months for recreation. The Tangipahoa River basin, including the research site, is located in a rural region dominated by agricultural use (dairy farms and livestock) with intermittent towns discharging treated effluent along the river. Large portions of the population in the watershed are not connected to municipal waste treatment systems. The site TR3 was chosen as the main sampling site because of its location downstream of most agricultural and domestic runoff/discharge. Additionally, a previous study of FC densities along the Tangipahoa River at sites TR1, TR3,



**Figure 1** | Locations of sample stations, Tangipahoa River, LA. Sample station TR3 corresponds to USGS hydrological station 07375500, located 30°30'23.19"N, 90°21'40.74"W.

TR4, TR5, TR6 and TR7 showed no significant differences between sites ( $p < 0.05$ ) (Gary W. Childers, unpublished data). This data was collected before 2005; therefore there may be some differences between modelled TR3 values and other site values. Faecal coliform densities have previously been reported to vary seasonally in this river (Anderson *et al.* 1990). Historically this river harbours densities of FC classifying it as unsuitable for primary contact recreational use (James 1987). In late 2004 – early 2005 treatment plants in several small municipalities were retrofitted with improved disinfection systems (Lake Pontchartrain Basin Foundation, personal communication).

### Sample collection and faecal coliform enumeration

A total of 884 FC samples were collected on 442 sampling days (replicate samples collected at the same time). Samples were collected on Monday mornings between approximately 05.00 h and 10.00 h at site TR3 from January 2000 to January 2008. Surface water grab samples were collected in duplicate sterile 1 l plastic containers and transported directly to the laboratory on ice for immediate analysis. Faecal coliform densities were enumerated using the most probable number (MPN) method according to *Standard Methods SM 9221 E (2005)* at Southeastern Louisiana University Microbiology Testing Laboratory (SLUMTL). Samples were serially diluted in phosphate buffered saline using a 10-fold dilution scheme to 10<sup>-6</sup>. A five tube MPN setup with A-1 medium was performed with duplicates for each sample and MPN values were averaged.

### NPMR of faecal coliform loading from storm water runoff

#### Selection of predictors for NPMR

A thorough treatment of the steps involved in creating a NPMR model is given in McCune (2006b). The Tangipahoa River has four real-time water stations (USGS) that monitor river stage, three that measure river stage and precipitation, and one station that measures river stage, precipitation and discharge (USGS station 07375500, TR3). All measurements are recorded twice per hour and available through the USGS National Water Information System web interface

(<http://waterdata.usgs.gov/nwis/rt>, accessed 23 November 2010). The predictor selection was limited to data obtainable through the USGS portal to increase the general applicability of the modelling procedure.

Daily discharge and river stage data from station TR3, and daily precipitation data from USGS stations 07375430 (TR7), 07375300 (TR10) and TR3 was retrieved from December 1999 to January 2008. Daily discharge data (cubic feet per second) was log transformed. No transformations were applied to precipitation data (inches day<sup>-1</sup>). A five day lag series (sum of antecedent rainfall or discharge up to five days) of the data was created to account for prior conditions and interactions associated with initial influx of storm water, and decreases in FC loading due to dilution after excessive precipitation/runoff (flushing).

### NPMR model search and evaluation

The predictor variables selected for model construction (Table 1) were expected to have interactions; however, the explicit form of the interactions was not known nor assumed *a priori*. Automatic interaction incorporation and model construction was accomplished by trial and error search. A NPMR local linear regression (LLR) with a Gaussian weighting kernel model search was implemented using Hyperniche<sup>®</sup> (MjM Software, Gleneden Beach, Oregon). Neighbourhood size is the number of observations incorporated into a local estimate; setting a minimum threshold helps ensure models

are not over-fitted to the data. The minimum neighbourhood size required for estimate was set to 5% of the total number of observations. For instance, a data point that has a unique combination of predictors not shared by at least 5% of the total number of observations would be omitted to help prevent erroneous conclusions based on limited information.

A leave-one-out cross-validation procedure was utilized to help eliminate over-fitting of models and provide an estimate of goodness of fit. The cross-validated  $R^2$  ( $\alpha R^2$ ) value was calculated from Equation (1), by eliminating point  $i$  from the estimate, for every data point used to construct the model.

$$\alpha R^2 = 1 - \frac{RSS}{TSS} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y}_i)^2} \quad (1)$$

The NPMR model fit ( $\alpha R^2$ ) was optimized using a kernel smoother in conjunction with a pre-specified local model (local linear) and varying the Gaussian weighting function. The tolerances defined the weighting function for each predictor and were equal to one SD of the Gaussian kernel smoother. Sensitivity analysis was performed by adjusting predictor values incrementally and observing the response. Both the predictor and responses were scaled using two methods of calculation (Sensitivity 1 and Sensitivity 2, details of calculations are described by McCune (2006b)) to produce a ratio where a value of 1 is equal to a 1:1 correspondence of predictor:response. The sensitivity provides an indication of how important a predictor is in the model. Additional statistical analyses of data were calculated using Excel (Microsoft Corporation, Redmond, WA) and Systat 11 (Systat Software, Chicago, IL).

## RESULTS AND DISCUSSION

### Long-term trends in faecal coliform densities at site TR3

A FC runoff model is only valid if FC loading is dominated by runoff-related factors rather than unpredictable point sources. Therefore, an important consideration for our

**Table 1** | Predictors for model search and forecasting

Predictor	Description
D	Mean daily discharge (Log10 cfs) at site TR3
D1	Mean daily discharge (Log10 cfs) at site TR3 day-1
D2	Mean daily discharge (Log10 cfs) at site TR3 day-2
D3	Mean daily discharge (Log10 cfs) at site TR3 day-3
D4	Mean daily discharge (Log10 cfs) at site TR3 day-4
P1	Sum of precipitation (inches day <sup>-1</sup> ) at sites TR3, TR7 and TR10 day-1
P2	Sum of precipitation (inches day <sup>-1</sup> ) at sites TR3, TR7 and TR10 day-2
P3	Sum of precipitation (inches day <sup>-1</sup> ) at sites TR3, TR7 and TR10 day-3
D <sub>est</sub>	Estimated daily discharge (Log10 cfs) at site TR3

model construction was the initial assessment of sources of faecal pollution entering the Tangipahoa River. The long-term trends of FC densities and modelled densities at site TR3 are presented in Figure 2. The observed FC rolling mean ( $n=25$ ,  $\text{Log}_{10}$  MPN  $100 \text{ ml}^{-1}$ ) varied at both annual and multi-year scales. FC loading displayed a cyclic pattern with peak mean densities occurring during winter months (high discharge months) and valleys occurring in summer months through the last four years in the time series.

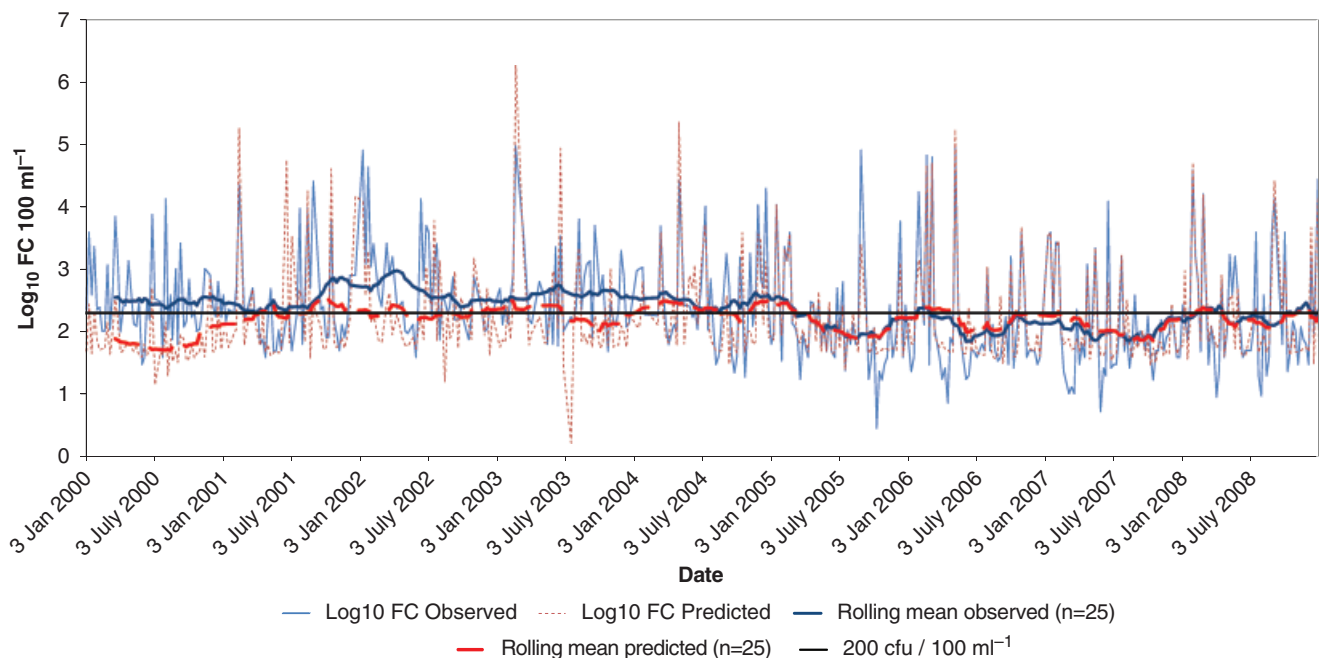
The annual seasonal trend was not readily apparent from 2000 to 2001, was more discernable from 2002 to 2004, and pronounced from 2005 to 2008. This result was not surprising for two reasons: 1) a drought occurred from 2000 to 2001; and 2) local sewage plants discharging treated effluent into the Tangipahoa River received upgraded disinfection systems during late 2004 (EPA 2008). Prior to 2005 the FC loading associated with storm water runoff may have been masked by point source discharges that contributed proportionately more FC during lower discharge events.

Lower FC values are clearly present in mean annual FC densities pre-2005 (Table 2). The FC concentration from 2000 to 2005 had a geometric mean of  $356 \text{ MPN } 100 \text{ ml}^{-1}$ , which

then significantly decreased to 139 after 2005 ( $p < 0.000$ ). The standard deviation of FC MPN  $\text{Log}_{10}$  CFU was lower in the years before 2005 (0.700 compared to 0.867), coinciding with the onset of a more pronounced seasonal pattern.

The percentage of samples that exceeded the primary recreational contact standard of  $200 \text{ FC MPN } 100 \text{ ml}^{-1}$  also declined from 56% to 29.6% pre- and post-January 2005, respectively. The minimum FC concentration observed before 2005 was  $19 \text{ MPN } 100 \text{ ml}^{-1}$ , after 2005 every year had a minimum value of less than  $9 \text{ MPN } 100 \text{ ml}^{-1}$ . The lowest value observed was  $3 \text{ MPN } 100 \text{ ml}^{-1}$  in 2005.

The long-term trends indicate that the point sources were of large enough magnitude during low river discharge periods (summer months) pre-2005 that input of FC from non-optimal disinfection and other point sources was sufficient to maintain a fairly consistent high FC load. During winter months, the increased discharge would have a diluting effect on faecal pollution if point sources were the only FC source; however, the increased FC loading due to surface runoff compensates for the point source dilution effect. As a result, we constructed two separate models: 1) a FC runoff model based on data from 2005 through 2008, excluding the



**Figure 2** | Time series of observed and modelled FC densities at station TR3. In 2005 there was a significant shift in observed values. The LLR-NPMR model was constructed using post-2005 data and accounts for the increased model agreement with post-2005 observations. The rolling mean values are centred averages that represent the previous 12 weeks of data, 12 weeks of data post the point, and the current point.

**Table 2** | Annual faecal coliform concentrations observed at site TR3

Year	No. samples	Log <sub>10</sub> MPN faecal coliform 100 ml <sup>-1</sup>				% samples above 200 FC 100 ml <sup>-1</sup>
		Avg	SD	Min	Max	
2000	49	2.505	0.601	1.477	4.123	59.2
2001	43	2.439	0.729	1.588	4.410	41.9
2002	41	2.703	0.738	1.588	4.906	65.9
2003	42	2.608	0.667	1.690	4.972	59.5
2004	43	2.515	0.770	1.276	4.419	53.5
2005	49	2.112	0.797	0.452	4.903	30.6
2006	47	2.115	0.974	0.858	4.952	25.5
2007	51	2.018	0.745	0.724	4.079	23.5
2008	49	2.332	0.934	0.954	4.588	38.8
2001–2005	218	2.551	0.700	1.276	4.972	56.0
2005–2008	196	2.143	0.867	0.452	4.952	29.6
Total	414	2.358	0.809	0.452	4.972	43.5

pre-2005 data; 2) a separate model for pre-2005 data. Even if some other factor (such as hydrologic differences pre-2005) was responsible for these observations, current conditions support using post-2005 data for model construction. These findings highlight the dependence of FC runoff models on low levels and small fluctuations in point-source pollution.

### Faecal coliform runoff model evaluation

The first LLR-NPMR model was constructed from post-2005 FC data using 196 sample dates. A best model was selected based on fit, predictor tolerances and sensitivities. The final model had an  $xR^2 = 0.57$  (Equation (1)) and is summarized in Table 3. The LLR-NPMR fitting procedure selected employed a Gaussian weighted window seeking to optimize  $xR^2$  based on multiplicative combination of predictors listed in Table 1. We also evaluated the model using the Nash-Sutcliffe coefficient of efficiency ( $E$ ) (Nash & Sutcliffe 1970), root mean square error (RMSE), mean absolute error (MAE), and index of agreement ( $d$ ) (Willmont 1981), all commonly employed to evaluate water quality models (Harmel & Smith 2007). To avoid under-estimation of error and over-estimation of goodness-of-fit, data observations that had a neighbourhood of less than five were excluded from traditional goodness of fit and error measurements. The LLR-NPMR had a RMSE = 0.43, a MAE = 0.31,  $E = 0.72$  and

$d = 0.91$ . The error estimates and other indices indicate that the LLR-NPMR model produced a good fit.

Tolerances represent one standard deviation of the Gaussian window used in the kernel smoother as estimated in the final model. It is the predictor values and observations combined with tolerances that define the LLR-NPMR

**Table 3** | Predictor tolerances and sensitivities, model  $xR^2 = 0.5721$ , avg neighbourhood size = 50.35

Predictor	Tolerance*	Min	Max	Sensitivity1 <sup>†</sup>	Sensitivity2 <sup>‡</sup>
D <sup>‡</sup>	0.5283	2.4742	3.9836	1.5801	1.9608
D <sub>1</sub>	0.2731	2.4728	3.8382	1.0618	1.5795
D <sub>2</sub>	1.081	2.4771	3.9186	0.9297	2.0162
D <sub>3</sub>	0.8527	2.4742	3.8954	1.0343	2.0376
D <sub>4</sub>	0.1386	2.4742	3.8603	0.9279	1.227
P <sub>1</sub> <sup>§</sup>	2.55	0	5.1	0.8576	2.2746
P <sub>2</sub>	4.492	0	5.99	0.4367	0.8656
P <sub>3</sub>	0.4785	0	9.57	0.9132	2.2592

\*Tolerance is defined as one SD of the Gaussian window utilized in producing the best estimate and is in units of the predictor.

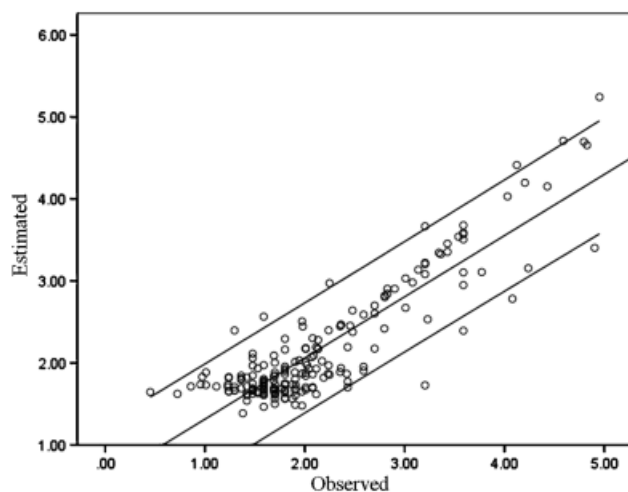
<sup>†</sup>Sensitivity is a measure of how much a predictor influences a response with a sensitivity of 1 equal to a 1:1 predictor to response ratio, sensitivity2 gives more weight to large differences (McCune 2006b).

<sup>‡</sup>D predictors are daily discharge (log cfs) at indicated day. D is the discharge at that day; D<sub>1</sub> is the previous day's discharge, etc.

<sup>§</sup>P predictors are summed precipitation data (inches day<sup>-1</sup>). P<sub>1</sub> is the sum of the previous day's precipitation; P<sub>2</sub> is the sum of two prior days' precipitation, etc.

model. The small tolerance of  $D$ ,  $D_1$ ,  $D_4$  and  $P_3$  coupled with strong to intermediate sensitivity indicate that these predictors have nonlinear global responses (Table 3). The presence of nonlinear global responses lends support for using a NPMR model with the data, as opposed to using multiple linear regression (MLR) procedures. The remaining predictors ( $D_2$ ,  $D_3$ ,  $P_1$  and  $P_2$ ) have larger tolerances with respect to observed ranges; however, their sensitivities were similar in magnitude to the other predictors, implying a more globally linear response for these variables.

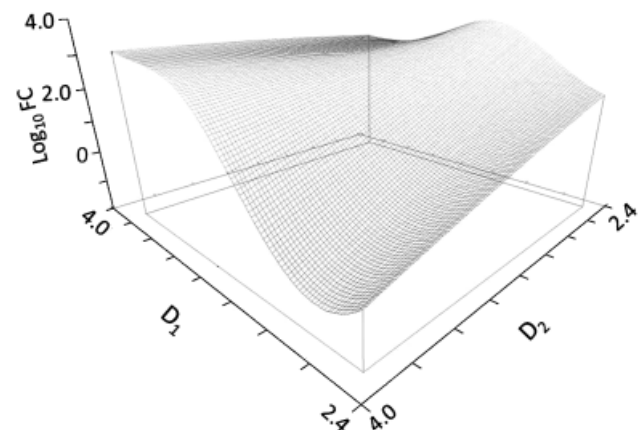
A plot of estimated vs. observed values (Figure 3) shows that the model has a larger amount of error associated with estimates that fall below 100 FC MPN  $100 \text{ ml}^{-1}$  and instances where the observed values greatly exceed predicted. From a public health and management decision viewpoint, the failure of the model to accurately predict FC values below 100 FC MPN  $100 \text{ ml}^{-1}$  is not significant, since these values are below regulatory limits. The lack of model fit during low flow/FC events may be more representative of variability not related to runoff during low flow periods, but rather due to point source variance. There is a cause for concern regarding the model's utility when either observed or predicted FC levels greatly exceed regulatory limits and are not in agreement. The observed values that greatly exceeded predicted values were investigated and it was found that in several cases an increase in rainfall occurred without an increase in



**Figure 3** | Plot of LLR-NPMR modelled data (estimated) and observed data. The lines represent 95% confidence intervals.

discharge; however, the corresponding FC levels increased by more than predicted by the model. In all cases where a rainfall spike without significant discharge spike occurred, the model predicted levels were still greater than regulatory limits ( $n=5$ , avg  $\text{Log}_{10}$  FC MPN observed = 4.20, and predicted = 2.93) with one exception. The case where the observed value ( $\text{Log}_{10}$  MPN = 3.20, predicted = 1.73) exceeded regulatory limits and the modelled value did not may be due to either insufficient rainfall data (rainfall did not occur at weather stations/rain gauge not properly functioning) or a point-source intrusion.

Another advantage of NPMR modelling is the ability to reconstruct response surfaces to evaluate predictor influence. Since responses are not always linear with respect to interactions, and the shape of the response is not always known beforehand, a tool to visualize and interpret interactions is important. Discharge events associated with storm water were not expected to produce a purely linear response with respect to FC levels, as storm events were expected to produce a peak in FC concentration followed by a dilution effect. To examine the utility of NPMR in producing response surfaces in a FC runoff model, a three-dimensional response surface was constructed (Figure 4) from previous day discharge ( $D_1$ ) and two day prior discharge ( $D_2$ ) to project the interaction between these variables with respect to FC levels. When  $D_2$  levels were high and  $D_1$  levels low, the lowest FC



**Figure 4** | Three-dimensional response surface of discharge interaction derived from the LLR-NPMR model. The response surface indicates how the previous day(s) discharge (one day previous discharge and two day prior discharge,  $D_1$  and  $D_2$ , respectively) interact and influence FC concentrations. For example, when  $D_2$  is high and  $D_1$  is low, then FC levels are low. This surface is not flat and therefore a linear model would have difficulty reproducing the model results of LLR-NPMR in certain scenarios.

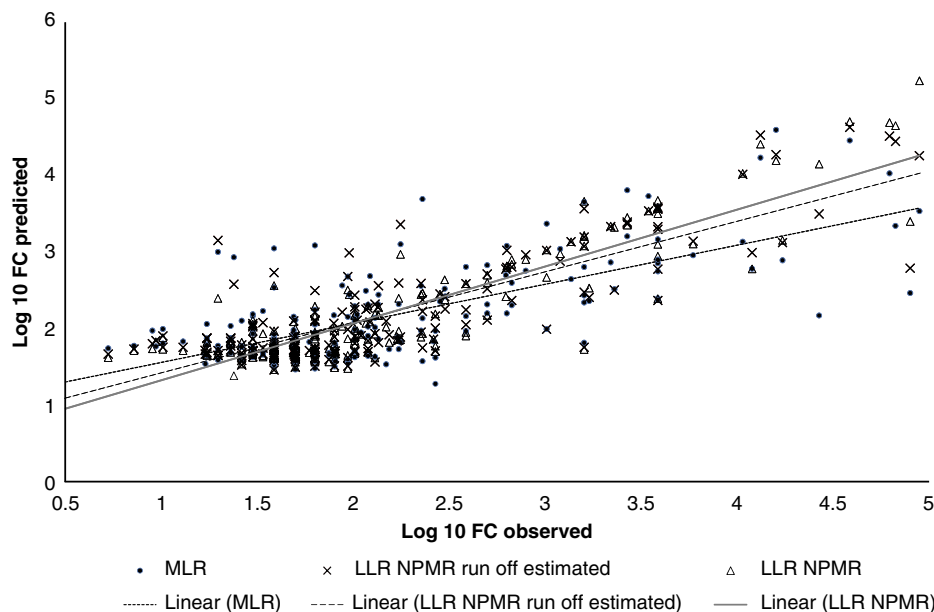
values were observed and, consequently, predicted. As expected,  $D_1$  had a greater influence on FC levels. For intermediate levels of  $D_1$  there was an increase in FC values that was inversely proportional to  $D_2$ . This illustrates and accounts for the casual observation that highest FC values were obtained when intermediate to high discharge occurred at  $D_1$ , and low discharge occurred at  $D_2$ . The response shape also emphasizes the likelihood that a linear model would inaccurately forecast FC levels, as a complex nonlinear function would be needed to reproduce such a response surface.

### Incorporation of discharge estimation into FC runoff model

The original LLR-NPMR model utilized actual daily discharge estimates that were available only post-event for FC estimation. This approach was used to give an optimal fit of data; however, the actual daily discharge is not available until the day after the event. A runoff model is more applicable when it can be utilized in real-time, such as when a recreational water body is in use. To address this, a discharge runoff model was constructed using the LLR NPMR approach in

order to estimate daily discharge in real-time. A total of 1,482 data discharge values (from USGS), in conjunction with precipitation data, was used to construct a discharge model to estimate  $D$  for a real-time FC run-off model. The hydrological model estimated  $D$  ( $D_{est}$ ) with a  $\chi R^2 = 0.93$  using all predictors in Table 1 except for  $D$ . Estimated discharge was then used to replace  $D$  in the original LLR-NPMR model. The resulting model was compared with a multivariate linear regression model (MLR), as well as the original LLR-NPMR model that utilized actual  $D$  values (Figure 5).

For a simple direct comparison of the different models,  $R^2$  values were calculated. The LLR-NPMR model using  $D_{est}$  performed better than the MLR model that utilized actual  $D$  values, but did not perform as well as the LLR-NPMR model that utilized the actual  $D$  value. Goodness-of-fit measures and error indices further confirmed the superior performance of the LLR-NPMR models (Figure 5). Furthermore, the MLR model had more instances where observed values exceeded recreational water quality standards, while predicted values did not. The use of  $D_{est}$  is fully justified by the results and suggests that its application may provide recreational users of the river with up-to-date estimates of water quality.



**Figure 5** | Comparison of multivariate linear regression (MLR) model with LLR-NPMR models (model with actual  $D$  value and model with  $D_{est}$ ). The root mean square error (RMSE), mean absolute error (MAE), Nash-Sutcliffe coefficient of efficiency ( $E$ ), and index of agreement ( $d$ ) was: 1) MLR 0.59, 0.44, 0.42 and 0.76; 2) LLR-NPMR 0.43, 0.31, 0.7 and 0.9; 3) LLR-NPMR with  $D_{est}$  0.5, 0.36, 0.59 and 0.85, respectively. All calculations were only done with estimates that had neighbourhood values  $> 5$  to avoid overestimation of model agreement with observations.

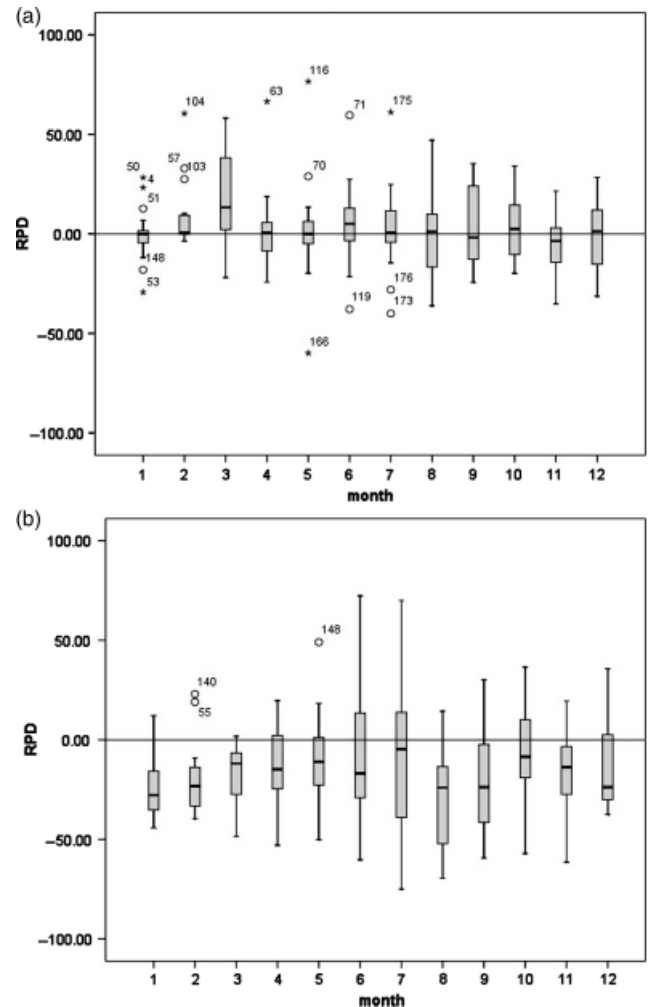


### Use of FC runoff model for long-term trend evaluation in a watershed

Although many of the long-term trends were identifiable by plotting the centred average (Figure 2), it is difficult to directly discern if these differences were due to external driving forces such as drought or point-source incursions. By back estimation of pre-2005 FC levels using the post-2005 LLR-NPMR, and vice-versa (forecasting of post-2005 FC using a similarly constructed NPMR model using pre-2005 data) it was determined that FC loading was greater than expected due to hydrologic predictors alone. Examination of model residuals revealed greater variability in residual values for pre-2005 compared with post-2005 ( $p < 0.000$ ). In contrast, the pre-2005 samples (actual observations) displayed less variance than the post-2005 samples (Table 2). These observations led us to conclude that a significant change in FC loading occurred in 2005, lowering the overall mean FC concentrations by reduction of point-source intrusions not solely influenced by rainfall and runoff. Owing to the seasonal trends observed in the FC data, it was important to assess the reliability of model predictions in relation to season. Auto-correlation plots of monthly means of relative percentage difference (RPD) of residuals showed no significant autocorrelation. This finding indicates that seasonal differences unrelated to discharge and precipitation were insignificant. Boxplots of RPD demonstrated that the model could accurately predict FC level regardless of season (Figure 6(a)). When the pre-2005 observations were compared with results of the post-2005 calibrated model, it was evident that the model consistently under-predicted FC concentrations (Figure 6(b)).

### CONCLUSIONS

The modelling approach adopted here for water quality monitoring offers numerous advantages compared with other methods. Namely, a reasonably representative model was created using a minimum of variables that are readily available via Internet resources. The method described here alleviates the need for complicated physical models, subsequent parameter estimation, and automatically accounts for interactions among predictors. While other nonparametric methods, such as ANNs, also allow for automatic interaction



**Figure 6** | Boxplots of monthly relative percentage difference (RPD) of residuals from the LLR-NPMR model using (a) post-2005 observations and (b) pre-2005 observations.

incorporation, how the interactions are derived is not readily apparent. The NPMR techniques allows for predictor interactions to be visualized. This study was a unique application of NPMR technique and the utility of this method should make it more commonplace in water quality modelling.

This modelling approach is adaptable to accommodate fairly complex interactions without the need for identifying and fitting all parameters as it is done automatically using NPMR. The results of this study highlight the importance of considering watershed characteristics before attempting to create FC (or other water quality parameter) runoff models. The most important consideration is whether the water body being studied is dominated by runoff FC loading or

point-source intrusions. This approach can also be used to monitor future changes in river water quality, as significant departures from modelled FC levels may be indicative of point-source incursions. The NPMR model framework is general enough to be applicable to a variety of watersheds and, as demonstrated in this study, can potentially estimate shifts in watershed characteristics associated with climate change and anthropogenic activities.

## ACKNOWLEDGEMENTS

This work was supported by funding from the USGS (Award No. 03HQAG0109) and the Southeastern Louisiana University Microbiology Testing Laboratory (SLUMTL). We would like to thank Dr Erin Watson-Horzelski for helpful comments on the manuscript. We would also like to acknowledge Andrea Bourgeois-Calvin and Ronnie Carter of the Lake Pontchartrain Basin Foundation for helpful discussions.

## REFERENCES

- Anderson, A. C., Abdelghani, A., Stumpf, J. & Rice, J. C. 1990 *Bacteriological criteria for recreational waters along the Tangipahoa River, DEQ Contract 24022-901-01*. Tulane University School of Public Health and Tropical Medicine, New Orleans.
- Arvanitidou, M., Papa, A., Constantinidis, T. C., Danielides, V. & Katsouyannopoulos, V. 1997 The occurrence of *Listeria* spp. and *Salmonella* spp. in surface waters. *Microbiol Res.* **152**(4), 395–397.
- Boehm, A., Ashbolt, N. J., Colford Jr, J. M., Dunbar, L. E., Fleming, L. E., Gold, M. A., Hansel, J. A., Hunter, P. R., Ichida, A. M., McGee, C. D., Soller, J. A. & Weisberg, S. B. 2009 *A sea change ahead for recreational water quality criteria*. *J. Water Health* **7**(1), 9–20.
- Brion, G. M. & Lingireddy, S. 1999 *A neural network approach to identify non-point sources of microbial contamination*. *Water Res.* **14**, 3099–3106.
- Chandramouli, V., Brion, G., Neelakantan, T. R. & Lingireddy, S. 2007 *Backfilling missing microbial concentrations in a riverine database using artificial neural networks*. *Water Res.* **41**, 217–227.
- Chandramouli, V., Neelakantan, T. R., Brion, G. & Lingireddy, S. 2008 *Predicting enteric virus presence in surface waters using artificial neural network models*. *Environ. Eng. Sci.* **25**, 53–62.
- Dufour, A. P. 1984 *Health effects criteria for fresh recreational waters. EPA-600/1-84-004*. US EPA, Washington, DC.
- Eleria, A. & Vogel, R. M. 2005 *Predicting fecal coliform bacteria levels in the Charles River, Massachusetts, USA*. *J. Am. Water Resour. Assoc.* **41**, 1195–1209.
- EPA 2008 *Nonpoint Source Program Success Story, Reducing Human and Animal Waste Discharge Restored Recreational Uses*. EPA 841-F-08-001Z. USEPA, Washington, DC.
- Harmel, R. D. & Smith, P. K. 2007 *Consideration of measurement uncertainty in the evaluation of goodness-of-fit in hydrologic and water quality modeling*. *J. Hydrol.* **337**, 326–336.
- Harwood, V. J., Levine, A. D., Scott, T. M., Chivukula, V., Lukasik, J., Farrah, S. R. & Rose, J. B. 2005 *Validity of the indicator organism paradigm for pathogen reduction in reclaimed water and public health protection*. *Appl. Environ. Microbiol.* **71**(6), 3163–3170.
- He, L. M. & He, Z. L. 2008 *Water quality prediction of marine recreational beaches receiving watershed baseflow and stormwater runoff in southern California, USA*. *Water Res.* **42**(10–11), 2563–2573.
- Horman, A., Rimhanen-Finne, R., Maunula, L., von Bonsdorff, C.-H., Torvela, N., Heikinheimo, A. & Hanninen, M.-L. 2004 *Campylobacter* spp., *Giardia* spp., *Cryptosporidium* spp., noroviruses, and indicator organisms in surface water in southwestern Finland, 2000–2001. *Appl. Environ. Microbiol.* **70**(1), 87–95.
- Hurst, C. J. & Murphy, P. A. 1996 *The transmission and prevention of infectious disease*. In: *Modeling Disease Transmission and its Prevention by Disinfection*, Hurst, C. J. (ed.), Cambridge University Press, Cambridge, UK, pp. 3–54.
- James, L. S. 1987 *An Analysis of the Sanitary Water Quality of the Tangipahoa River and a Survey of Selected Tangipahoa Streams*. Southeastern Louisiana University, Hammond, LA.
- Jamieson, R., Gordon, R., Joy, D. & Lee, H. 2004 *Assessing microbial pollution of rural surface waters: A review of current watershed scale modeling approaches*. *Agr. Water Manage.* **70**, 1–17.
- Karine Lemarchand, P. L. 2003 *Occurrence of Salmonella spp. and Cryptosporidium spp. in a French coastal watershed: relationship with fecal indicators*. *FEMS Microbiol. Lett.* **218**(1), 203–209.
- Kashefipour, S. M., Lin, B., Harris, E. & Falconer, R. A. 2002 *Hydro-environmental modelling for bathing water compliance of an estuarine basin*. *Water Res.* **36**(7), 1854–1868.
- Leclerc, H., Mossel, D. A. A., Edberg, S. C. & Struijk, C. B. 2001 *Advances in the bacteriology of the coliform group: their suitability as markers of microbial water safety*. *Annu. Rev. Microbiol.* **55**(1), 201–234.
- Lin, B., Kashefipour, S. M. & Falconer, R. A. 2003 *Predicting near-shore coliform bacteria concentrations using ANNS*. *Water Sci Technol.* **48**(10), 225–232.
- Mahloch, J. L. 1974 *Comparative analysis of modeling techniques for coliform organisms in streams*. *Appl. Microbiol.* **27**, 340–345.
- Mas, D. M. L. & Ahlfeld, D. P. 2007 *Comparing artificial neural networks and regression models for predicting faecal coliform concentrations*. *Hydrolog. Sci.* **52**, 713–731.
- McCarthy, D. T., Mitchell, V. G., Deletic, A. & Diaper, C. 2007 *Escherichia coli in urban stormwater: explaining their variability*. *Water Sci. Technol.* **56**(11), 27–34.
- McCune, B. 2006a *Non-parametric habitat models with automatic interactions*. *J. Veg. Sci.* **17**, 819–830.
- McCune, B. 2006b *Nonparametric Multiplicative Regression for Habitat Modeling*, <http://www.pcord.com/NPMRintro.pdf> (accessed 22 November 2010).
- Nash, J. E. & Sutcliffe, J. V. 1970 *River flow forecasting through conceptual models, Part I: a discussion of principals*. *J. Hydrol.* **10**, 282–290.
- Payment, P., Berte, A., Prevost, M., Menard, B. & Barbeau, B. 2000 *Occurrence of pathogenic microorganisms in the Saint Lawrence*

- River (Canada) and comparison of health risks for populations using it as their source of drinking water. *Can. J. Microbiol.* **46**(6), 565–576.
- Prüss, A. 1998 Review of epidemiological studies on health effects from exposure to recreational water. *Int. J. Epidemiol.* **27**(1), 1–9.
- Setti, I., Rodriguez-Castro, A., Pata, M. P., Cadarso-Suarez, C., Yacoubi, B., Bensmael, L., Moukrim, A. & Martinez-Urtaza, J. 2009 Characteristics and dynamics of *Salmonella* contamination along the Coast of Agadir, Morocco. *Appl. Environ. Microbiol.* **75**(24), 7700–7709.
- Standard Methods for the Examination of Water and Wastewater* 2005 21st edition, American Public Health Association/American Water Works Association/Water Environment Federation, Washington, DC.
- Steets, B. M. & Holden, P. A. 2003 A mechanistic model of runoff-associated fecal coliform fate and transport through a coastal lagoon. *Water Res.* **37**, 589–608.
- Wade, T. J., Pai, N., Eisenberg, J. N. & Colford, J. M., Jr 2003 Do US Environmental Protection Agency water quality guidelines for recreational waters prevent gastrointestinal illness? A systematic review and meta-analysis. *Environ. Health Perspect.* **111**(8), 1102–1109.
- Whitman, R. L., Nevers, M. B., Korinek, G. C. & Byappanahalli, M. N. 2004 Solar and temporal effects on *Escherichia coli* concentration at a Lake Michigan swimming beach. *Appl. Environ. Microbiol.*, **70**(7), 4276–4285.
- Willmont, C. J. 1981 On the validation of models. *Phys. Geogr.* **2**, 184–194.

First received 16 February 2010; accepted in revised form 26 August 2010. Available online 6 January 2011