# MARLBOROUGH DISTRICT COUNCIL

# Predictive Models for Swimming Sites in Marlborough

January 2025

# Predictive Models for Swimming Sites in Marlborough

MDC Technical Report No. 25-002

File Reference/Record No: E370-007-001/255649

January 2025

Report Prepared by:

**Steffi Henkel**

Senior Environmental Scientist - Water Quality

Environmental Science & Monitoring Group

&

**Nandakumar Thachapilly**

Environmental Data Analyst

Environmental Science & Monitoring Group

Approved by

**Alan Johnson**

Environmental Science & Monitoring Manager

# Executive Summary

During the warmer months, Marlborough District Council monitors the most popular swimming spots weekly to assess the risk to swimmers from waterborne diseases. The results are made available on the LAWA website, but water quality is highly variable, and the data is not suitable to provide users with up-to-date information on the health risk at a particular site. To overcome this, predictive models can be used and the LAWA website provides the ability to display model outputs.

Analysis of the data showed that the datasets are highly skewed, with most data points in the low range of indicator bacteria concentrations and comparatively few data points indicating unsafe swimming conditions. The unbalanced data set limits the model approaches that can be used. However, simplifying the predictions into two categories ("safe" and "unsafe") allows the use of logistic regression or boosted tree analysis despite the data limitations.

Four model approaches were tested: Random Forest Classifier, Stochastic Gradient Boost, XGBoost, and Logistic Regression. Of these, Logistic Regression showed the greatest promise and was subsequently used.

For all but two sites, rainfall or river flow are generally the best predictors of indicator bacteria concentrations. The exceptions were the monitoring sites on the lower Taylor River, which are heavily influenced by localized sources such as wildfowl and potential sewage/stormwater cross contamination.

Three different preceding rainfall totals (12hr, 24hr, 48hr) from several rainfall stations near the swimming sites, as well as river flows (for river sites), were used as predictors for swimming safety.

Table 1 shows the rainfall sites and associated rainfall totals (column 2) that best predicted if a site was safe for swimming. The right side of the table shows rainfall cut-offs for different probabilities. A probability of 0.5 means there is a 50% chance of unsafe bacteria levels.

| | Swimming Site | Predictor Variable | Rainfall totals (in mm) for different probabilities of unsafe bacteria levels | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.8 |
| Coastal | Anakiwa | Kaituna Rainfall - 12 hr | 24.8 | 33.8 | 39.8 | 44.7 | 49.2 | 64.6 |
| | Momorangi Bay | Waikawa Rainfall - 24 hr | 22.8 | 49.8 | 67.8 | 82.5 | 96.0 | 142.2 |
| | Ngakuta Bay | Kaituna Rainfall - 12 hr | 7.7 | 12.4 | 15.6 | 18.2 | 20.6 | 28.7 |
| | Picton Foreshore | Kaituna Rainfall - 24 hr | 9.4 | 18.4 | 24.4 | 29.3 | 33.8 | 49.2 |
| | Waikawa Bay | Kaituna Rainfall - 12 hr | 18.0 | 25.4 | 30.3 | 34.3 | 38.0 | 50.6 |
| | Waikutakuta/Robin Hood Bay | Rarangi Rainfall - 24 hr | 16.4 | 21.8 | 25.4 | 28.4 | 31.1 | 40.3 |
| | Pukatea/Whites Bay | Rarangi Rainfall - 24 hr | 21.4 | 25.7 | 28.5 | 30.9 | 33.0 | 40.3 |
| Rivers | Rai River at Rai Falls | Rai Rainfall - 24 hr | 10.8 | 18.2 | 23.1 | 27.1 | 30.8 | 43.4 |
| | Te Hoiere/Pelorus Rv at Totara Flat | Rai Rainfall - 24 hr | 20.6 | 26.9 | 31.0 | 34.4 | 37.5 | 48.2 |
| | Te Hoiere/Pelorus Rv at Pelorus Bridge | Tunakino Rainfall - 24 hr | 88.8 | 115.8 | 133.8 | 148.5 | 162.0 | 208.2 |
| | Wairau Rv at Ferry Bridge | Blenheim Rainfall - 48 hr | 16.0 | 19.7 | 22.1 | 24.2 | 26.0 | 32.3 |
| | Wairau Rv at Blenheim Rowing Club | Blenheim Rainfall - 24 hr | 10.1 | 13.2 | 15.3 | 17.0 | 18.5 | 23.9 |
| | Waihopai Rv at Craighlochart #2 | Spray Rainfall - 24 hr | 7.1 | 9.7 | 11.4 | 12.7 | 14.0 | 18.3 |

**Table 1: Rainfall totals for different probability cut-offs, based on the best performing models for the individual monitoring sites.**

Choosing a cut-off involves balancing minimizing risk to swimmers (lower probabilities) and reducing instances of predicting unsafe conditions when they are actually safe (higher probabilities). As the predictions concern human health, it is sensible to choose a lower probability for unsafe bacteria concentrations, such as a probability of 0.2.

Regular monitoring of the sites should continue to ensure no significant short-term or long-term changes are occurring. Using the additional data collected each season, fine-tuning the existing model and further exploring the use of alternative approaches with different adjustment methods during model training and validation is necessary to keep the model outputs current. This should be done regularly, for example following every summer season. Additional monitoring during rainfall should be conducted whenever possible to provide better data for future model refinements.

# Table of Contents

## Table of Figures

## List of Tables

# 1.    Introduction

Recreational activities such as swimming and boating in rivers and coastal areas are an essential part of the Kiwi lifestyle, especially during the summer months. While water quality is generally good, there are times when swimmers may be at risk from waterborne diseases, which can cause symptoms such as stomach aches and diarrhea.

To assess these risks, the council monitors the most popular swimming spots during the warmer months of the year. Monitoring involves weekly sampling from the beginning of November until the end of March. These samples are analysed for indicator bacteria: E. coli for rivers and Enterococci for coastal waters. The concentrations of these bacteria are compared against national guidelines, which are designed to indicate the risk of contracting illnesses such as Campylobacteriosis, the most commonly reported waterborne disease in New Zealand. The guidelines aim to protect about 95-99% of users, though some individuals may be more susceptible than others.

The guideline document [10] provides threshold values for each type of indicator bacteria. Based on these values, sample results are categorized into three "Modes." Concentrations within the "Green Mode" indicate a low health risk to swimmers. If bacteria levels reach the "Alert Mode," the infection risk increases slightly. Although swimming is still considered safe, this is a signal for the council to investigate potential sources of faecal pollution. Once bacteria concentrations exceed the "Action Mode" threshold, the health risk is deemed unacceptable. Table **2** shows the ranges of indicator bacteria concentrations for each "Mode".

| Mode | Coastal | Rivers | Meaning |
|---|---|---|---|
| | Enterococci/ 100mL | E. coli/ 100mL | |
| **Green (Safe) Mode** | <140 | <260 | **Safe** for contact recreation |
| **Amber (Alert) Mode** | 140 - 280 | 260 - 550 | **Increased risk** for health, but still considered safe |
| **Red (Action) Mode** | >280 | >550 | **Unsafe** for contact recreation |

**Table 2**: **Modes in the Microbiological Water Quality Guidelines for Marine and Freshwater Recreational Areas [10].**

Once water samples have been taken, they are sent to Christchurch for analysed by an independent lab (Hill Labs), and results are received the day after sampling. Once available, the results can be viewed on the LAWA.org.nz website. Many users rely on the most recent results as an indicator of the current health risk to swimmers, assuming the data represents conditions until the next samples are taken. However, this is not an appropriate use of the information. Bacteria concentrations can fluctuate rapidly, sometimes within hours or even minutes. Therefore, the sample results only represent conditions at the specific date and time the sample was collected.

Instead, users are encouraged to refer to the general site gradings (Suitability for Contact Recreation Grades – SFR Grades) to inform their decisions. These grades are based on data collected over several summer seasons and provide a broader assessment of a site's suitability for swimming. More information about the SFR Grades can be found in the council's monitoring reports [7] and the national guideline document [10].

Many swimming sites have been monitored for over 15 years, and the data shows that high-risk bacteria levels usually occur after rainfall or during high river flows. As a result, the Marlborough

District Council and the local Health Board recommend avoiding swimming for at least 48 hours after rainfall or when water appears visibly turbid.

Although this recommendation is shared with the public at the beginning of each summer season through media releases and information on the LAWA.org.nz website, many swimmers remain unaware. Additionally, visitors to swimming sites may not know whether it has rained recently in the area they are visiting, and finding rainfall information can be challenging.

Fortunately, the relationship between rainfall or river flow and bacteria concentrations can be used to produce real-time predictions of health risks through modelling. While models always involve some uncertainty, careful selection of parameters and the use of a precautionary approach when setting thresholds can provide predictions that greatly improve upon the information currently available to users.

# 2. Model Choice

There are various models available for predicting the concentration of indicator bacteria, including linear regression, Bayesian networks, and machine learning approaches [5; 6]. Most of these models aim to predict bacteria concentrations, which are then used to assess the risk to swimmers.

These modeling techniques typically require a wide range of sample results across different bacteria concentrations, along with associated predictor variables, and a relatively large number of data points [16]. However, because water quality at most swimming sites in Marlborough is generally good, bacteria concentrations are often low, with many samples showing levels below the detection limit (see Figure 1).



**Figure 1**: **Examples of the distribution of data (histograms) across indicator bacteria concentrations. The example sites are those with the best and worst recreational water quality for the two different types of sites (rivers and coastal beaches).**

As Figure 1 shows, that generally the number of sample results within the Action Mode (unsafe for swimming) is comparatively small even for sites with comparatively poor grading. This creates a challenge, as individual results at higher concentrations can disproportionately influence the model. To avoid this, a better approach is to categorize the data rather than using the actual concentration values. The guidelines provide three categories, but these can be simplified into two: concentrations

that are considered safe for swimming (Green and Amber Modes) and those that indicate an unacceptable health risk (Red Mode).

By categorizing the data in this way, alternative modeling techniques, such as logistic regression or classification trees, can be employed. Logistic regression produces a probability function and is often used in medical research to predict the likelihood of illness or mortality based on factors such as age, exposure, and habits [2,13], but it has also been used in predictive models for recreational water quality [6].

Classification trees are commonly used for prediction of recreational water quality [6, 17] and are able to also provide probability outputs. We trialed three different classification tree model approaches, Random Forest Classifier, Stochastic Gradient Boost and XGBoost.

Despite catorization of the data, the unbalance in the dataset remains and can still lead to loss of the information in the category with few datapoints (the "unsafe" category) and overfitting of the category with the majority of the datapoints (the "safe'' category). This needs to be carefully managed during model training.

# 3.    Methodology

E. coli and Enterococci concentration results were categorized as either "0" for values in the Green and Amber Modes, or "1" for values in the Red (unsafe) Mode. Data from limited time periods when higher indicator bacteria levels were observed, but not linked to rainfall or river flow, were excluded. Examples are high Enterococci concentrations in Momorangi Bay due to a wastewater pipe breakage, or elevated Enterococci levels in Waikutakuta/Robin Hood Bay caused by large accumulations of vegetation along the shore.

Where possible, data spanning 15 years was used. Although longer datasets were available for most sites, significant trends over time impacted model performance. However, even within the shorter 15-year period, trends were observed at some sites. Except for the Waihopai River at Craiglochart #2, these trends were of relatively small magnitude [7] and could therefore be ignored. For the Waihopai River, trend analysis indicated an annual increase of 6 n/100mL, resulting in an overall rise of 90 n/100mL over the 15-year period. To adjust for this trend, a linear adjustment equivalent to the annual change was applied to the data. Although the actual change may not have been linear and reflects median rather than individual values, this adjustment provides a reasonable approximation in the absence of more specific data.

Predictive models were developed for all long-term monitoring sites in the recreational water quality program, except for the Taylor River. The Taylor River, an urban river swimming site, is influenced by additional factors beyond rainfall, such as sewage infrastructure failures, wildfowl movements, and dog droppings along the river reserve.

The models were developed as single-predictor-variable models. The predictor variables tested were:

- Total rainfall at nearby monitoring stations within the 12, 24, and 48 hours prior to sampling for indicator bacteria.
- River flow at or near the swimming site at the time of sampling for indicator bacteria.

Model fit was evaluated using ROC curves and the associated area under the curve (AUC)[1], as well as p-values calculated using the Wald test. All analyses was conducted in R using the caret, glmnet and pROC libraries.

---

[1] A ROC (Receiver Operating Characteristic) curve provides a visual representation of a model's performance across all data points. The AUC is the area under the ROC curve and higher values generally indicate better model performance.

Using a subset of sites, the performance of the four model approaches[2] was tested. Overall, logistic regression provided the most robust models, and further refinement was based on outputs from logistic models only. The best fitting model for each site was chosen using the already mentioned model performance parameters as well as predicted vs. observed curves and confusion matrixes. Figure 2 shows an example of the performance measures used to assess model performance.
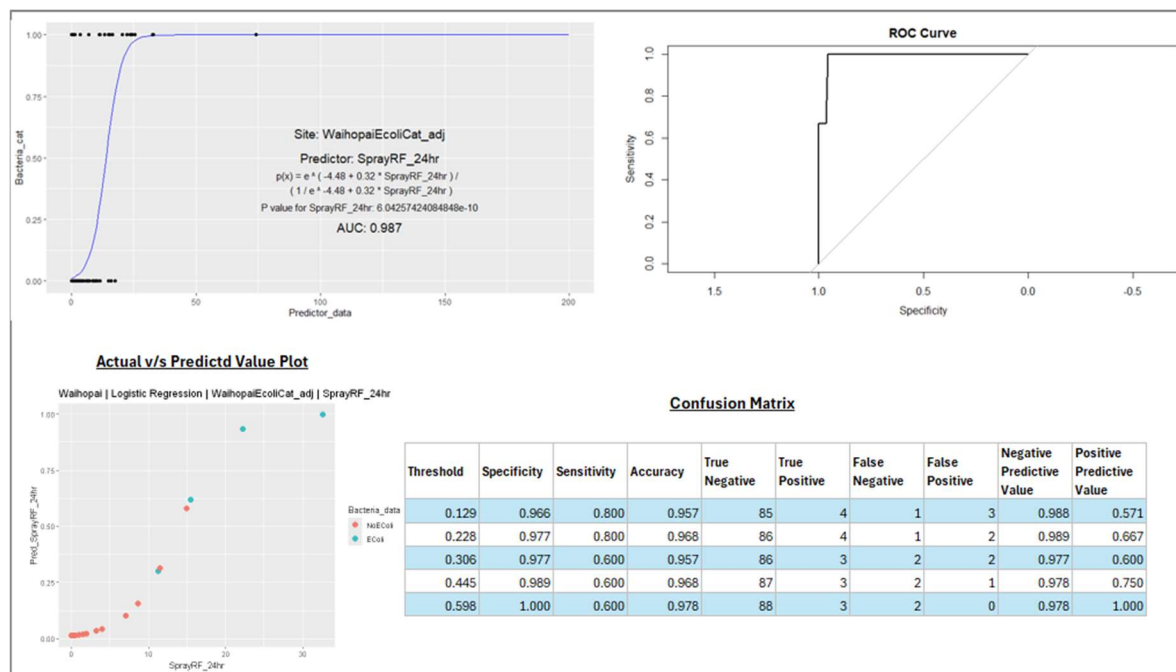


**Figure 2: Example of performance parameters and graphical representations used to determine the best fitting models for each site.**

# 4.     Model Results and Discussion

The model parameters and performance measures [4, 8] for the best-fitting models for each monitoring site are shown in Table 3 (parameters and performance measures for all models can be found in Appendix 6.1).

| | Swimming Site | Predictor Variable | Intercept | Predictor Estimate | P-value | AUC |
|---|---|---|---|---|---|---|
| Coastal | Anakiwa | Kaituna Rainfall - 12 hr | -4.43 | 0.09 | 0.011 | 0.706 |
| | Momorangi Bay | Waikawa Rainfall - 24 hr | -2.88 | 0.03 | 0.007 | 0.891 |
| | Ngakuta Bay | Kaituna Rainfall - 12 hr | -3.50 | 0.17 | <0.001 | 0.842 |
| | Picton Foreshore | Kaituna Rainfall - 24 hr | -3.04 | 0.09 | <0.001 | 0.808 |
| | Waikawa Bay | Kaituna Rainfall - 12 hr | -4.18 | 0.11 | <0.001 | 0.946 |
| | Waikutakuta/Robin Hood Bay | Rarangi Rainfall - 24 hr | -4.66 | 0.15 | <0.001 | 0.993 |
| | Pukatea/Whites Bay | Rarangi Rainfall - 24 hr | -6.27 | 0.19 | <0.001 | 0.998 |
| Rivers | Rai River at Rai Falls | Rai Rainfall - 24 hr | -3.39 | 0.11 | <0.001 | 0.929 |
| | Te Hoiere/Pelorus Rv at Totara Flat | Rai Rainfall - 24 hr | -4.88 | 0.13 | <0.001 | 0.991 |
| | Te Hoiere/Pelorus Rv at Pelorus Bridge | Tunakino Rainfall - 24 hr | -4.86 | 0.03 | 0.005 | 0.953 |
| | Wairau Rv at Ferry Bridge | Blenheim Rainfall - 48 hr | -5.72 | 0.22 | <0.001 | 0.993 |
| | Wairau Rv at Blenheim Rowing Club | Blenheim Rainfall - 24 hr | -4.82 | 0.26 | <0.001 | 0.976 |
| | Waihopai Rv at Craighlochart #2 | Spray Rainfall - 24 hr | -4.48 | 0.32 | <0.001 | 0.987 |

**Table 3: Covariates and Model parameters for the best performing model for each monitoring sites**

---

[2] Random Forest Classifier, Stochastic Gradient Boost, XGBoost and Logistic regression.

For the two Wairau River monitoring sites, river flow also resulted in a comparatively good model fit. However, because the LAWA website mainly supports rainfall data models, the rainfall models were chosen for these sites.

The poorest-performing model was for Anakiwa, located at the head of Queen Charlotte Sound/Tōtaranui near Okiwa Bay Estuary. The predominance of fine sediment in this shallow area potentially acts as a reservoir for indicator bacteria [14, 18], a factor not accounted for in the current rainfall-based model. Incorporating tidal influences [3] could improve model performance and should be considered in future developments.

While the LAWA website uses 24-hour rainfall data as input, for many sites the 12-hour rainfall statistic proved to be a better predictor of unsafe bacteria concentrations. It is suggested that the feasibility of using 12-hour rainfall on the LAWA website be explored, potentially through the creation of virtual measurements in Hilltop.

| | Swimming Site | Predictor Variable | Rainfall totals (in mm) for different probabilities of unsafe bacteria levels | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.8 |
| Coastal | Anakiwa | Kaituna Rainfall - 12 hr | 24.8 | 33.8 | 39.8 | 44.7 | 49.2 | 64.6 |
| | Momorangi Bay | Waikawa Rainfall - 24 hr | 22.8 | 49.8 | 67.8 | 82.5 | 96.0 | 142.2 |
| | Ngakuta Bay | Kaituna Rainfall - 12 hr | 7.7 | 12.4 | 15.6 | 18.2 | 20.6 | 28.7 |
| | Picton Foreshore | Kaituna Rainfall - 24 hr | 9.4 | 18.4 | 24.4 | 29.3 | 33.8 | 49.2 |
| | Waikawa Bay | Kaituna Rainfall - 12 hr | 18.0 | 25.4 | 30.3 | 34.3 | 38.0 | 50.6 |
| | Waikutakuta/Robin Hood Bay | Rarangi Rainfall - 24 hr | 16.4 | 21.8 | 25.4 | 28.4 | 31.1 | 40.3 |
| | Pukatea/Whites Bay | Rarangi Rainfall - 24 hr | 21.4 | 25.7 | 28.5 | 30.9 | 33.0 | 40.3 |
| Rivers | Rai River at Rai Falls | Rai Rainfall - 24 hr | 10.8 | 18.2 | 23.1 | 27.1 | 30.8 | 43.4 |
| | Te Hoiere/Pelorus Rv at Totara Flat | Rai Rainfall - 24 hr | 20.6 | 26.9 | 31.0 | 34.4 | 37.5 | 48.2 |
| | Te Hoiere/Pelorus Rv at Pelorus Bridge | Tunakino Rainfall - 24 hr | 88.8 | 115.8 | 133.8 | 148.5 | 162.0 | 208.2 |
| | Wairau Rv at Ferry Bridge | Blenheim Rainfall - 48 hr | 16.0 | 19.7 | 22.1 | 24.2 | 26.0 | 32.3 |
| | Wairau Rv at Blenheim Rowing Club | Blenheim Rainfall - 24 hr | 10.1 | 13.2 | 15.3 | 17.0 | 18.5 | 23.9 |
| | Waihopai Rv at Craighlochart #2 | Spray Rainfall - 24 hr | 7.1 | 9.7 | 11.4 | 12.7 | 14.0 | 18.3 |

**Table 4: Rainfall totals for different probability cut-offs. A common cut-off is at the probability of 0.5, which equates to a 50:50 chance of indicator bacteria levels showing unsafe swimming conditions. A more precautionary cut-off (i.e, 0.2, which equates to a 20% chance of unsafe conditions) could be chosen to protect swimmers.**

Logistic regression produces a probability curve, which estimates the likelihood of unsafe bacteria levels based on predictor variables (e.g., 24-hour rainfall). A common cut-off for unsafe conditions is a 0.5 probability, meaning a 50% chance of unsafe bacteria concentrations. However, different cut-offs can be selected depending on priorities. Choosing a higher probability cut-off reduces false positives (where conditions are predicted unsafe when they are safe). Conversely, a lower probability cut-off ensures more instances of unsafe conditions are captured, at the cost of increasing false positives. If the priority is protecting swimmers' health, a lower cut-off, such as 0.2 (indicating a 20% chance of unsafe bacteria), might be selected. Table 4 shows rainfall thresholds for various probability cut-offs for the selected models.

Once an appropriate cut-off is chosen, the corresponding rainfall total can be used to inform the LAWA model. The model will mark a site as unsafe once rainfall exceeds this threshold. Additionally, the LAWA module applies the rainfall limit over the preceding 48 hours, ensuring that sites remain marked as unsafe for at least 48 hours, aligning with the recommendation to avoid swimming for 48 hours after rainfall.

Figure 3 show an example of the difference in the information that would be displayed on LAWA if model outputs were to be used. The example uses the information from the Te Hoiere/Pelorus at

Totara Flat from the previous summer season. During that season, LAWA showed the site to be unsafe for swimming during one week in March 2024. As the display was based on sampling results, the information was updated only once per week. Bacteria concentrations in the samples taken during that season are shown as black bars in the top graph. For the only sample with high bacteria levels (in March 2024), the status update is delayed as sample analysis takes at least 24 hours. It can also be seen that the weekly sampling often misses rainfall events (rainfall is shown in blue), which would cause high bacteria level unsafe for swimmers. For example, the rainfall event in early December, was one that would cause unsafe swimming conditions with high certainty, but was only capture through sampling at the tail end of the event when bacteria levels were already decreasing again.

The information that would have been displayed on LAWA if predictive models were used is shown in the middle graph. Predictive models would significantly improve the information provide to the users of the website by capturing all events that are caused by rainfall run-off. That is despite the chance of occasionally showing a site to be unsafe for swimming when it might not be, such as the rainfall event in January 2024.
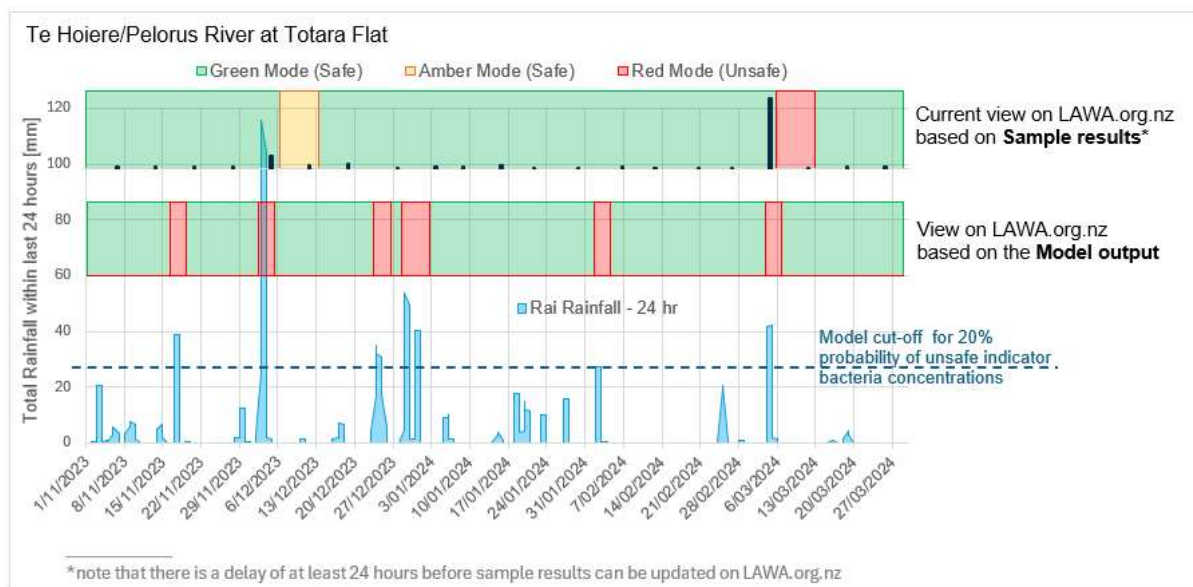


**Figure 3: Example of the difference in information displayed on the LAWA website based on the latest sample results (current display) and the use of predictive model outputs (proposed display)**

Although indicator bacteria concentrations can now be predicted using these models, ongoing site monitoring remains essential [6]. This is not only mandated by the National Policy Statement for Freshwater Management but is also crucial for detecting changes over time, whether due to long-term water quality improvements or temporary issues such as sewage infrastructure malfunctions.

If additional funding and resources become available, targeted sampling during rainfall events would enhance model accuracy. Such data could support the use of more advanced models, such as Bayesian networks [1] and artificial neural networks [5, 12], which have the potential to yield more precise predictions. Regardless of additional sampling, regular model reviews and updates using newly collected data from routine monitoring are necessary. Additionally, model adjustment methods to deal with the unbalanced datasets should be further investigated to provide more representative model outputs [9, 15].

# 5.   References

1.   Avila R et al. (2018) Evaluating statistical model performance in water quality prediction. Journal of Environmental Management. 2018 Jan 15:206:910-919. DOI: 10.1016/j.jenvman.2017.11.049

2.   Boateng EY & Abaye DA (2019) A Review of the Logistic Regression Model with Emphasis on Medical Research. Journal of Data Analysis and Information Processing, Vol.7 No.4,

3.   Boehm AB & Weisberg SB (2005) Tidal forcing of enterococci at marine recreational beaches at fortnightly and semidiurnal frequencies. Environmental Science & Technology 39 (15): 5575-5583.

4.   Castro, HM & Ferreira JC (2022) Linear and logistic regression models: when to use and how to interpret them? Jornal Brasileiro de Pneumologia 48(6).

5.   Diane MLM & Ahlfeld DP (2007) Comparing artificial neural networks and regression models for predicting faecal coliform concentrations. Hydrological Sciences Journal 52 (4): 713-731.

6.   Heasley C et al (2021) Systematic review of predictive models of microbial water quality at freshwater recreational beaches. PLoS One. 2021; 16(8): e0256785

7.   Henkel S (2023) Recreational Water Quality Report 2022-2023. Marlborough District Council Technical Report 23-005.

8.   James G et al (2021) An Introduction to Statistical Learning with Applications in R. Springer

9.   Krawcyk B (2016) Learning from imbalanced data: open challenges and future directions. Springer. DOI 10.1007/s13748-016-0094-0

10.   MfE/MoH (2003) Microbiological Water Quality Guidelines for Marine and Freshwater Recreational Areas. Ministry for the Environment.

11.   Raner M (2020) On logistic regression and a medical application. U.U.D.M. Project Report 2020:40.

12.   Rustam F et al. (2022) An Artificial Neural Network Model for Water Quality and Water Consumption Prediction. Water 2022, 14(21), 3359.

13.   Schober P & Vetter TR (2021) Logistic Regression in Medical Research. Anaesthesia & Analgesia 132(2):365-366. DOI:10.1213/ANE.0000000000005247

14.   Solo-Gabriele HM et al. (2016) Beach sand and the potential for infectious disease transmission: observations and recommendations. Journal of the Marine Association of the United Kingdom 96(1), 101-120.

15.   Van den Goorberg R et al. (2022) The harm of class imbalance corrections for risk prediction models: illustration and simulation using logistic regression. Journal of the American Medical Informatics Association, 29(9): 1525–1534

16. Van der Ploeg T et al. (2014) Modern modelling techniques are data hungry: a simulation study for predicting dichotomous endpoints. BMC Medical Research Methodology 14: 137.

17. Xu T et al. (2020) A predictive model of recreational water quality based on adaptive synthetic sampling algorithms and machine learning. Water Research Vol. 177.

18. Yamahara KM et al. (2007). Beach sands along the California coast are diffuse sources of faecal bacteria to coastal waters. Environmental Science & Technology 41 (13): 4515-4521.

# 6.    Appendix

## 6.1.    All Model Results

The Table below shows the model parameters for all site/covariate pairs tested, with the chosen models highlighted.

**Table 5:  Model parameters for all site/covariate pairs tested**

| Site | Covariate | Intercept | Predictor Estimate | P-value | AUC |
|---|---|---|---|---|---|
| Anakiwa_Cat | WaikawaRF_12hr | -4.31 | 0.06 | 0.025 | 0.698 |
| Anakiwa_Cat | WaikawaRF_24hr | -4.21 | 0.03 | 0.396 | 0.646 |
| Anakiwa_Cat | WaikawaRF_48hr | -4.19 | 0.02 | 0.650 | 0.612 |
| Anakiwa_Cat | KaitunaRF_12hr | -4.43 | 0.09 | 0.011 | 0.706 |
| Anakiwa_Cat | KaitunaRF_24hr | -4.51 | 0.06 | 0.008 | 0.654 |
| Anakiwa_Cat | KaitunaRF_48hr | -4.51 | 0.04 | 0.061 | 0.622 |
| Momorangi_Cat | WaikawaRF_12hr | -2.91 | 0.06 | 0.004 | 0.835 |
| Momorangi_Cat | WaikawaRF_24hr | -2.88 | 0.03 | 0.068 | 0.891 |
| Momorangi_Cat | WaikawaRF_48hr | -2.85 | 0.02 | 0.296 | 0.820 |
| Momorangi_Cat | KaitunaRF_12hr | -2.87 | 0.07 | 0.029 | 0.796 |
| Momorangi_Cat | KaitunaRF_24hr | -2.87 | 0.03 | 0.118 | 0.853 |
| Momorangi_Cat | KaitunaRF_48hr | -2.87 | 0.02 | 0.349 | 0.799 |
| NgakutaBay_Cat | WaikawaRF_12hr | -3.41 | 0.11 | <0.001 | 0.730 |
| NgakutaBay_Cat | WaikawaRF_24hr | -3.45 | 0.05 | <0.001 | 0.695 |
| NgakutaBay_Cat | WaikawaRF_48hr | -3.6 | 0.04 | <0.001 | 0.651 |
| NgakutaBay_Cat | KaitunaRF_12hr | -3.5 | 0.17 | <0.001 | 0.842 |
| NgakutaBay_Cat | KaitunaRF_24hr | -3.6 | 0.08 | <0.001 | 0.813 |
| NgakutaBay_Cat | KaitunaRF_48hr | -3.74 | 0.05 | <0.001 | 0.769 |
| PictonFors_Cat | WaikawaRF_12hr | -2.73 | 0.07 | 0.002 | 0.784 |
| PictonFors_Cat | WaikawaRF_24hr | -2.81 | 0.04 | 0.001 | 0.789 |
| PictonFors_Cat | WaikawaRF_48hr | -2.9 | 0.03 | 0.001 | 0.788 |
| PictonFors_Cat | KaitunaRF_12hr | -2.91 | 0.16 | <0.001 | 0.801 |
| PictonFors_Cat | KaitunaRF_24hr | -3.04 | 0.09 | <0.001 | 0.808 |
| PictonFors_Cat | KaitunaRF_48hr | -2.97 | 0.04 | <0.001 | 0.745 |
| Waikawa_Cat | WaikawaRF_12hr | -4.45 | 0.13 | <0.001 | 0.797 |
| Waikawa_Cat | WaikawaRF_24hr | -4.14 | 0.06 | 0.001 | 0.770 |
| Waikawa_Cat | WaikawaRF_48hr | -4.16 | 0.04 | 0.013 | 0.760 |
| Waikawa_Cat | KaitunaRF_12hr | -4.18 | 0.11 | <0.001 | 0.946 |
| Waikawa_Cat | KaitunaRF_24hr | -4.44 | 0.1 | <0.001 | 0.913 |
| Waikawa_Cat | KaitunaRF_48hr | -4.32 | 0.04 | <0.001 | 0.873 |
| RobinH_Cat | RarangiRF_12hr | -4.7 | 0.22 | <0.001 | 0.991 |
| RobinH_Cat | RarangiRF_24hr | -4.66 | 0.15 | <0.001 | 0.993 |
| RobinH_Cat | WaikawaRF_12hr | -4.83 | 0.24 | <0.001 | 0.985 |
| RobinH_Cat | WaikawaRF_24hr | -4.7 | 0.14 | 0.001 | 0.993 |
| WhitesBay_Cat | RarangiRF_12hr | -6.49 | 0.25 | <0.001 | 0.990 |
| WhitesBay_Cat | RarangiRF_24hr | -6.27 | 0.19 | <0.001 | 0.988 |
| WhitesBay_Cat | WaikawaRF_12hr | -7.14 | 0.26 | 0.006 | 0.988 |
| WhitesBay_Cat | WaikawaRF_24hr | -6.82 | 0.18 | 0.002 | 0.976 |

| Site | Covariate | Intercept | Predictor Estimate | P-value | AUC |
|---|---|---|---|---|---|
| RAR1_Cat | Tunak_12hr | -2.84 | 0.13 | <0.001 | 0.670 |
| RAR1_Cat | Tunak_24hr | -3.32 | 0.09 | <0.001 | 0.875 |
| RAR1_Cat | Tuank_48hr | -3.45 | 0.05 | <0.001 | 0.894 |
| RAR1_Cat | RaiRF_12hr | -2.93 | 0.15 | <0.001 | 0.717 |
| RAR1_Cat | RaiRF_24hr | -3.39 | 0.11 | <0.001 | 0.929 |
| RAR1_Cat | RaiRF_48hr | -3.53 | 0.06 | <0.001 | 0.927 |
| RAR1_Cat | Rai_Flow | -3.15 | 0.06 | <0.001 | 0.779 |
| RAR1_Cat | PelorTotara_Flow | -3.41 | 0.03 | <0.001 | 0.812 |
| RAR1_Cat | PelorBryants_Flow | -3.47 | 0.05 | <0.001 | 0.833 |
| TotaraF_Cat | Tunak_12hr | -3.68 | 0.16 | <0.001 | 0.957 |
| TotaraF_Cat | Tunak_24hr | -5.27 | 0.12 | <0.001 | 0.987 |
| TotaraF_Cat | Tuank_48hr | -4.3 | 0.05 | <0.001 | 0.947 |
| TotaraF_Cat | RaiRF_12hr | -3.84 | 0.17 | <0.001 | 0.973 |
| TotaraF_Cat | RaiRF_24hr | -4.88 | 0.13 | <0.001 | 0.991 |
| TotaraF_Cat | RaiRF_48hr | -4.25 | 0.05 | <0.001 | 0.962 |
| TotaraF_Cat | Rai_Flow | -3.8 | 0.07 | <0.001 | 0.891 |
| TotaraF_Cat | PelorTotara_Flow | -4.39 | 0.04 | <0.001 | 0.974 |
| TotaraF_Cat | PelorBryants_Flow | -4.59 | 0.06 | <0.001 | 0.981 |
| PelorusBr_Cat | Tunak_12hr | -6.27 | 0.13 | 0.006 | 0.940 |
| PelorusBr_Cat | Tunak_24hr | -4.86 | 0.03 | 0.005 | 0.953 |
| PelorusBr_Cat | Tuank_48hr | -4.77 | 0.02 | 0.090 | 0.869 |
| PelorusBr_Cat | RaiRF_12hr | -5.42 | 0.1 | 0.006 | 0.951 |
| PelorusBr_Cat | RaiRF_24hr | -4.8 | 0.03 | 0.006 | 0.922 |
| PelorusBr_Cat | RaiRF_48hr | -4.71 | 0.02 | 0.091 | 0.873 |
| PelorusBr_Cat | Rai_Flow | -4.32 | 0.01 | 6.627 | 0.673 |
| PelorusBr_Cat | PelorTotara_Flow | -4.46 | 0.01 | 1.032 | 0.682 |
| PelorusBr_Cat | PelorBryants_Flow | -4.56 | 0.01 | 0.262 | 0.692 |
| FerryBr_Cat | BlhmRF_12hr | -4.09 | 0.27 | <0.001 | 0.625 |
| FerryBr_Cat | BlhmRF_24hr | -4.99 | 0.32 | <0.001 | 0.812 |
| FerryBr_Cat | BlhmRF_48hr | -5.72 | 0.22 | <0.001 | 0.993 |
| FerryBr_Cat | WairauFlow | -5.4 | 0.01 | <0.001 | 0.965 |
| BlenheimRC_Cat | BlhmRF_12hr | -4.23 | 0.23 | 0.001 | 0.965 |
| BlenheimRC_Cat | BlhmRF_24hr | -4.82 | 0.26 | <0.001 | 0.976 |
| BlenheimRC_Cat | BlhmRF_48hr | -5.62 | 0.19 | <0.001 | 0.942 |
| BlenheimRC_Cat | WairauFlow | -5.22 | 0.01 | <0.001 | 0.981 |
| WaihopaiCat_adj | CraiglRF_12hr | -2.88 | 0.16 | 0.032 | 0.612 |
| WaihopaiCat_adj | CraiglRF_24hr | -3.47 | 0.23 | <0.001 | 0.784 |
| WaihopaiCat_adj | CraiglRF_48hr | -3.34 | 0.09 | <0.001 | 0.864 |
| WaihopaiCat_adj | SprayRF_12hr | -3 | 0.22 | 0.002 | 0.623 |
| WaihopaiCat_adj | SprayRF_24hr | -4.48 | 0.32 | <0.001 | 0.987 |
| WaihopaiCat_adj | SprayRF_48hr | -3.65 | 0.09 | <0.001 | 0.923 |
| WaihopaiCat_adj | Waihopai_Flow | -3.74 | 0.06 | <0.001 | 0.780 |