

## Prediction of Red Tide Events in the Philippines using MODIS-derived Oceanographic Parameters and XGBoost

Christian Angelo P. Arellano<sup>1</sup>, Althea Andrienne C. Castillo<sup>1</sup>, Ayin M. Tamondong<sup>1</sup>, Jommer M. Medina<sup>1</sup>, John Emmanuel D. Escoto<sup>1</sup>

<sup>1</sup>Department of Geodetic Engineering, University of the Philippines Diliman, Philippines - (christianangeloarellano, aandriennec@gmail.com, amtamondong, jmmedina, jdescoto)@up.edu.ph

**Keywords:** Fisheries Management Areas (FMA), prediction model, remote sensing.

### Abstract

Red tide, a commonly used misnomer for HABs in the Philippines, poses a significant threat to the environment, fisheries, and public health. Since current red tide detection methods, such as in-situ sampling, are mostly reactive and usually result in delayed issuance of advisories, this study developed a model for predicting red tide occurrences in the Philippines using XGBoost and MODIS-derived oceanographic parameters. Five key parameters, namely Chlorophyll-a (chl-a), Sea Surface Temperature (SST), Photosynthetically Available Radiation (PAR), Diffuse Attenuation Coefficient ( $K_d(490)$ ), and Particulate Backscattering Coefficient ( $b_{bp}(443)$ ), were extracted from MODIS Aqua 8-day composite products spanning 2003 to 2021. These were integrated with historical data from BFAR, covering the same period, to train a predictive model using the XGBoost algorithm. The final model demonstrated moderate performance as reflected in its accuracy (58%), F1-score (59%), and AUC (61%), with chl-a and  $K_d(490)$  as the most influential features based on their SHAP values. Its precision and recall of 58-59% showed its balanced predictive ability across classes, namely, banned and lifted. Model performance across different FMAs and seasons varied due to factors, such as minor variation in parameter values across adjacent FMAs and seasons, missing pixel values of crucial parameters, mismatched parameter values and red tide-linked conditions, and unevenly distributed training data. Among all, the model produced the most reliable and representative results in FMA 7, and the poorest in FMAs 10 and 11.

### 1. Introduction

#### 1.1 Background of the Study

As the primary food source for aquatic organisms, phytoplankton is a vital component of marine ecosystems. However, when natural and anthropogenic factors, including sewage discharge, agricultural activities, and surface runoff, disrupt the equilibrium in these ecosystems, the nutrient levels in marine ecosystems become excessively high, resulting in the development of Harmful Algal Blooms (HABs). Algal blooms, whether toxic or non-toxic, degrade water bodies and cause extensive fish kills, which can result from either toxin production or oxygen depletion (Pal et al., 2020; Khan et al., 2021).

In the Philippines, red tide is a commonly used misnomer to refer to HABs. However, based on the scientific definition, HABs refer to any algal proliferation that produces toxins or harmful effects on humans and marine animals, while red tide pertains to excessive phytoplankton growth, particularly dinoflagellates that produce saxitoxin, leading to water discoloration (Wexler, 2014; Azanza et al., 2024). In the Philippines, the most commonly reported toxic manifestation of red tide is paralytic shellfish poisoning (PSP), which can be acquired from consuming shellfish contaminated with *Pyrodinium bahamense*, *Gymnodinium catenatum*, or toxic *Alexandrium spp* (Azanza et al., 2024). From 1983 to 2002, the country recorded the highest number of PSP cases in Asia, with 2,124 incidents and 120 fatalities (Ching et al., 2015). Beyond health impacts, red tide outbreaks have also posed serious economic consequences.

Although there are already studies on detecting red tide, there remains a gap in research on predicting red tide events. Currently, the Bureau of Fisheries and Aquatic Resources (BFAR) heavily relies on water samples and shellfish meat sample collections to issue advisory bulletins on red tide (Yñiguez et al., 2020). However, sampling and analysis typically happen after requests

from local government units or after residents observe or experience the effects of red tide, making these approaches inefficient. Moreover, the diverse and fragmented coastal geography of the Philippines makes it difficult to model red tide events using ground-based data alone, as it requires extensive sampling to capture spatial variations. In contrast, remotely sensed datasets provide broader and consistent coverage for quantifying these differences across various coastal environments. Therefore, in this study, five (5) oceanographic parameters derived from MODIS imagery and historical data from BFAR were used by the XGBoost algorithm to predict red tide occurrences across the Philippines.

#### 1.2 Objectives

Given that existing methods of detecting red tide are mostly reactive, there is a need for predictive approaches to mitigate adverse effects on the environment, local fisheries, and public health. Thus, this research aims to develop a prediction model for red tide in the Philippines using MODIS-derived oceanographic parameters and XGBoost. The specific objectives of the study are:

1. To analyze the relationship between red tide events and five (5) MODIS-derived parameters, including Chlorophyll-a (chl-a), Sea Surface Temperature (SST), Photosynthetically Available Radiation (PAR), Diffuse Attenuation Coefficient ( $K_d(490)$ ), and Particulate Backscattering Coefficient ( $b_{bp}(443)$ ) using Pearson correlation and Seasonal and Trend decomposition using LOESS (STL).
2. To develop and assess the predictive performance of the XGBoost model for occurrences of red tide by comparing its predictions with the test set using performance metrics, such as precision, accuracy, recall, F1-score, and AUC.
3. To evaluate spatial and temporal patterns of predicted red tide bans and potential model biases across FMAs

and seasons using MODIS-derived oceanographic parameters and XGBoost.

organisms, toxin levels, and cell densities. However, not all these details were consistently available for every year.

### 1.3 Significance of the Study

This study is critical in addressing the need for predicting red tide events in the Philippines as part of a proactive approach to red tide monitoring and management. This approach not only addresses the gap in the prediction of red tide events but also offers a foundational step toward developing a proactive, scalable, and efficient tool for protecting communities and fisheries in the Philippines. Moreover, this approach aligns with various Sustainable Development Goals (SDGs) that can benefit the country: SDG 14 (Life Below Water), Targets 14.1 and 14.2, by reducing marine pollution and promoting sustainable marine resource management; SDG 3 (Good Health and Well-Being), Target 3.9, by mitigating health risks associated with red tide-related contamination; SDG 1 (No Poverty), Target 1.5, by preserving economic livelihoods and reducing vulnerability to the impacts of red tide; and SDG 8 (Decent Work and Economic Growth), Target 8.5, by fostering productive employment in the fishing and aquaculture sectors.

## 2. Methodology

### 2.1 General Workflow

Figure 1 outlines the general processes involved in this study, with each step described in the subsequent sections.

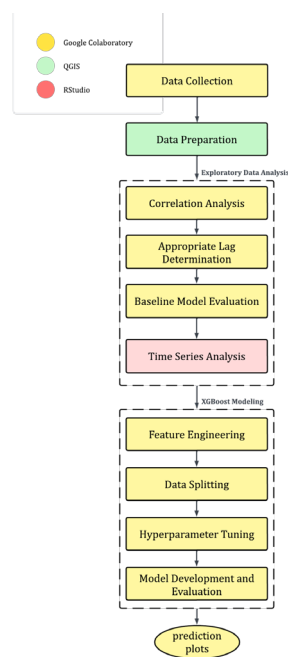


Figure 1. General Flowchart of the Study

### 2.2 Data Collection

To predict red tide events in the Philippines, various datasets were utilized that capture both environmental conditions and spatial distributions relevant to red tide occurrences. Historical records were obtained from the Bureau of Fisheries and Aquatic Resources - Fisheries Resources Management Division (BFAR-FRMD), based on bulletin advisories from 2003 to 2021. These included specific location and status, and in some cases, other relevant information for certain years, including causative

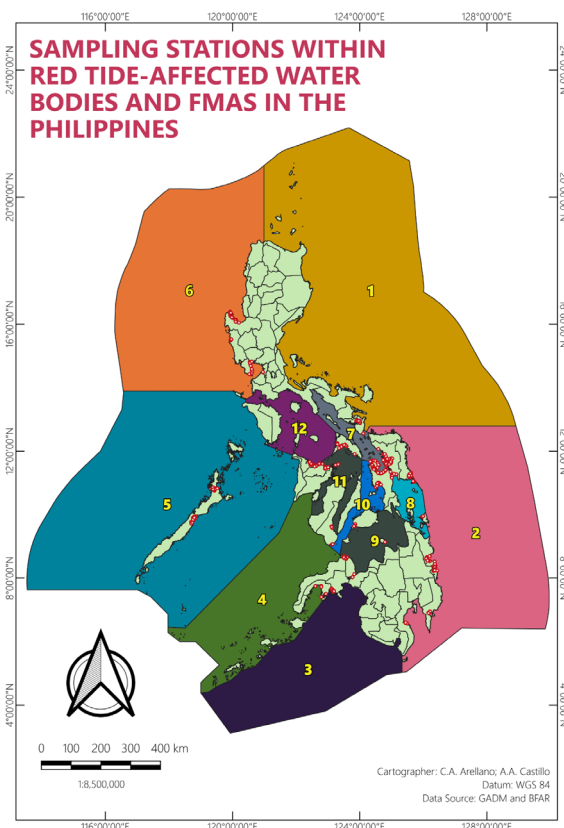


Figure 2. Sampling Stations of the Recorded Areas for Red Tide and Fisheries Management Areas (FMAs) in the Philippines (2002-2023)

Since spatial variability is crucial in this study, the researchers incorporated the coordinates of the 382 sampling stations of the 62 recorded areas in the bulletins over the 19 years, as seen in Figure 2. However, the bulletins did not specify which sampling station detected individual red tide events. To assess the environmental conditions associated with red tide events, Level 3 Standard Mapped Image (SMI) products were used for five (5) oceanographic parameters, specifically chl-a, which indicates phytoplankton biomass; SST, which affects phytoplankton growth;  $K_d(490)$ , which measures light reduction with depth;  $b_{bp}(443)$ , which measures light reflected by suspended particles; and PAR, which estimates sunlight for ocean photosynthesis. These 8-day composite datasets are satellite-derived and averaged on a consistent spatial grid. A total of 873 images were used for each parameter, except for SST, with 850 images, covering 2003 to 2021. They are also already land and cloud masked since Level 2 data, from which Level 3 products are derived, undergo quality assessment. In this way, it is ensured that Level 3 data contain only valid geophysical measurements, making it reliable for water-related studies such as red tide prediction. However, given its 4.6 km spatial resolution, mismatches with the actual location of sampling stations may occur, prompting analysis at the red tide-affected site level. For spatial visualization, additional geospatial layers were included, such as administrative boundaries from GADM, FMA boundaries from BFAR, to allow red tide-affected sites to be grouped based on similar oceanographic conditions, and the delineation of red tide-affected bays, coastal waters, and sampling stations.

## 2.3 Preparation of the Remotely-Sensed Input Parameters

**2.3.1 Reprojection:** All NASA Ocean Color Level 3 8-day composite NetCDF files were initially loaded into QGIS 3.40.3 to identify and extract the relevant subdatasets for each parameter. Then, a script was used to automate the bulk loading of usable subdatasets into grouped layers by parameter in QGIS. To ensure spatial consistency, all layers were reprojected to WGS 84/UTM Zone 51N. In this way, the downloaded datasets will be overlaid properly and expressed in consistent measurement units for subsequent processes and analysis.

**2.3.2 Extraction of Site-Specific Environmental Data:** Pixel values for the five (5) oceanographic parameters were extracted across all available dates for each parameter. This was done to obtain site-specific environmental conditions relevant to red tide events based on the pixel values corresponding to each sampling station using their coordinates. From here on, the extracted pixel values for each parameter were saved separately in CSV format for subsequent analysis.

## 2.4 Exploratory Data Analysis (EDA)

**2.4.1 Correlation Analysis:** In this study, pairwise Pearson correlation was used to analyze relationships between oceanographic parameters. The Pearson correlation coefficient, represented as  $r$ , for two variables,  $x$  and  $y$ , with paired values, is given by Equation 1:

$$r = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum(x_i - \bar{x})^2 \sum(y_i - \bar{y})^2}} \quad (1)$$

where  $x_i, y_i$  = individual data points  
 $\bar{x}, \bar{y}$  = means of each variable

This analysis helped identify patterns and dependencies between variables. For example, since chlorophyll-a (chl-a) is closely related to HAB distribution (Izadi et al., 2021), any parameter closely correlated with chl-a can also be a good proxy of red tide occurrence. This analysis also supported feature selection by pointing out the most significant parameters and confirmed preliminary hypotheses about drivers of red tide.

**2.4.2 Appropriate Lag Determination:** To determine the lag period for predictive model training, the historical records of red tide events in 2003 and the five (5) parameters were examined. Using Google Colaboratory, time series plots were produced for each bay and coastal water using the corresponding values of each parameter in 2003. The given banned and lifted classifications, as well as their corresponding dates of sample collection, were overlaid on the time series plots. Based on the plots, the lag period for each parameter was estimated. For chl-a, this was measured as the delay between the sample collection date and the first ban advisory. For sites with multiple sample collection dates, the first ban after the earlier collection period was used as the reference point for the next estimation. The same logic was applied to other parameters except SST. For SST, the lag period was measured by identifying when it reached the optimal range of 28°C to 30°C for the growth of *P. bahamense*, the most common causative organism (Smayda, 1997; Folio & Yap-Dejeto, 2022). The gap was calculated from the point of optimal SST value until the first ban advisory after the chl-a collection date, since it is the most commonly used proxy for algal blooms (Manzar Abbas et

al., 2019; Izadi et al., 2021; Joshi et al., 2024). This delay was a critical factor in estimating the lag period for each parameter because it affected the timing of red tide alerts.

**2.4.3 Baseline Model Evaluation:** To gauge the potential model performance of XGBoost using the given datasets, Lazy Predict was utilized. Lazy Predict is a library in Python that facilitates model training and evaluation of model performance. In this study, it was used to gain insight into the baseline values of evaluation metrics that could be achieved by XGBoost given the available data without any tuning process. The performance metrics tested include accuracy, balanced accuracy, ROC AUC, F1 score, and runtime. The results helped set realistic expectations and served as a benchmark for further improvements, particularly for XGBoost.

**2.4.4 Time Series Analysis:** Time series analysis serves as a valuable preliminary method for examining the trend and seasonality in datasets with chronologically arranged observations (Donatelli et al., 2022). To conduct the analysis, the extracted pixel values for each site-parameter combination were plotted and exported as individual images through RStudio. This involved generating time series plots of each parameter and decomposing the data into trend, seasonality, and residual components. Afterward, the plots were visually analyzed based on the trend and seasonality of each parameter.

More specifically, the Seasonal and Trend decomposition using LOESS (STL) method was used in R version 4.4.2 to break down time series objects into trend, seasonal, and remainder components. The seasonal component shows the recurring patterns that happen at regular intervals; the trend component provides the overall direction of the data; and the remainder component captures the random variations or residuals (Gordan et al., 2024). To extract these components, a custom R version 4.4.2 script was used to apply STL to all five parameters, and they were plotted and saved. The median value of all 8-day composite pixel values in a month was used to represent values for each month per parameter-site combination, addressing missing pixel values in some oceanographic parameters, which could be due to cloud cover, as STL cannot handle trailing missing values. Thus, a frequency of 12, representing 12 months, was used. In instances where an 8-day period overlapped two months, it was assigned to the month it was temporally closest to.

## 2.5 XGBoost Modeling

This study leveraged the loss and regularization capabilities of XGBoost, iteratively minimizing residual errors and preventing overfitting. Its capability to handle missing data and proven success in past red tide prediction studies further support its selection as the algorithm for this study (Izadi et al., 2021).

**2.5.1 Feature Engineering:** Since the parameters do not have the same range of values, feature scaling, which is a technique employed to standardize the data, was implemented using mean normalization, as shown by Equation 2. This ensured that model training is more efficient. To capture temporal patterns and historical trends, lag features were created using the determined 1-period lag, corresponding to the previous 8-day value for each parameter.

$$x' = \frac{x - \mu}{\max(x) - \min(x)} \quad (2)$$

where  $x$  = individual value of a feature  
 $\mu$  = average mean of the values of a feature  
 $\min(x)$  = minimum feature value

$\max(x)$  = maximum feature value

**2.5.2 Data Splitting:** Using a time-based stratified data splitting method, the 2003 to 2021 extracted pixel values for each 8-day composite oceanographic parameter and their corresponding 1-period lag values for every unique site, were split into training (80%), validation (10%), and testing (10%) datasets, the most common ratio for data split and helps prevent data leakage (Delisi, 2024). Given the complex dynamics of red tide events, allocating a larger portion for training is ideal to allow the model to learn intricate patterns in the data better. To avoid bias and ensure reliable performance, a balanced representation of banned and lifted classifications was preserved across all three datasets.

**2.5.3 Hyperparameter Tuning:** To enhance the predictive performance of the model, hyperparameter tuning was done before every model training. This was performed using Optuna in Google Colaboratory, which employed Random Search. A hundred trials using this method were performed, selecting random hyperparameter values from a set range rather than going through all possible hyperparameter combinations. To prevent possible overfitting, the tuning process was also configured to stop early when the F1-score did not show improvement over the trials.

**2.5.4 Model Development and Evaluation:** In each trial, after sampling a specific hyperparameter combination, the model was trained using the mean values and lagged values of Chl-a, SST, PAR,  $K_d(490)$ , and  $b_{bp}(443)$  per site from 2003 to early 2020. Model performance was evaluated using AUC-ROC to guide hyperparameter selection. The best combination was then used to retrain the model and generate prediction probabilities for red tide classifications ("1" for banned, "0" for lifted) in the validation set covering late 2020 to early 2021. These probabilities were converted into binary outcomes using the optimal threshold from Youden's J Statistic. Predictions were marked "TRUE" or "FALSE" based on actual classifications. Model performance was assessed through a confusion matrix (Table 1), using accuracy (Equation 3), recall (Equation 4), precision (Equation 5), F1-score (Equation 6), macro average (Equation 7), and weighted average (Equation 8) as metrics. The macro average equally weights all classes, while the weighted average accounts for class imbalance. It was also assessed against two benchmarks: (1) Inokuchi et al. (2024), whose red tide classification model achieved 52–58% across metrics; and (2) Lazy Predict results, representing baseline performance without tuning. SHAP (SHapley Additive exPlanations) was also used to interpret feature contributions. Final evaluation on the test set, covering the unseen data from 2020–2021, used the same metrics and SHAP, with results visualized across FMAs.

	Classified as banned	Classified as lifted
Banned	True Positive (TP)	False Negative (FN)
Lifted	False Positive (FP)	True Negative (TN)

Table 1. Confusion Matrix

$$Accuracy = \frac{TP + TN}{TP + FN + FP + TN} \quad (3)$$

$$Recall = \frac{TP}{TP + FN} \quad (4)$$

$$Precision = \frac{TP}{TP + FP} \quad (5)$$

$$F1 - Score = \frac{2 \left( \frac{TP}{TP + FP} \right) \left( \frac{TP}{TP + FN} \right)}{\left( \frac{TP}{TP + FP} \right) + \left( \frac{TP}{TP + FN} \right)} \quad (6)$$

$$Macro\ Average = \frac{1}{N} \sum_{i=1}^N Metric_i \quad (7)$$

$$Weighted\ Average = \frac{\sum_{i=1}^N (Metric_i \times Support_i)}{\sum_{i=1}^N Support_i} \quad (8)$$

where

$N$  = total number of classes

$Metric_i$  = performance metric for class  $i$

$Support_i$  = number of samples in class  $i$

### 3. Results and Discussions

#### 3.1 Exploratory Data Analysis (EDA)

##### 3.1.1 Relationship Among the Oceanographic Parameters:

The strong positive relationship between chl-a and  $K_d(490)$ , as observed in Figure 3, suggests that when the phytoplankton biomass gets high, it inhibits light penetration in the water column, which has been noted in other studies. Moderate correlation between chl-a and  $b_{bp}(443)$  implies that when chl-a levels increase, the concentration of suspended particles also increases, resulting in increased light backscattering. This also highlights that  $b_{bp}(\lambda)$  is sensitive to both phytoplankton and non-algal particles, limiting the strength of the correlation.

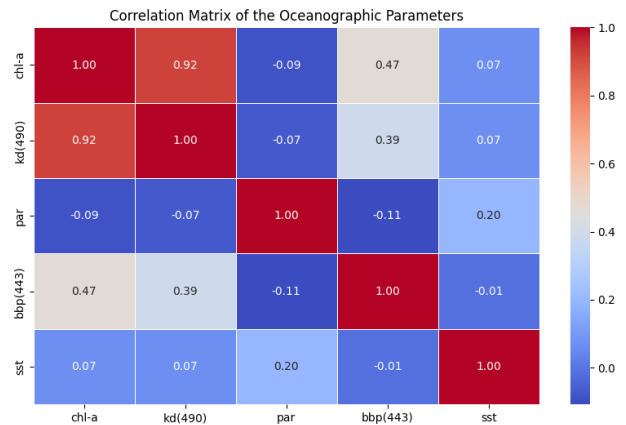


Figure 3. Correlation Matrix of All Oceanographic Parameters

Contrary to the results, a weak negative correlation between chl-a and SST has been observed, following the overall trend in some equatorial and tropical oceans where higher SST is normally linked with higher stratification, which lowers nutrient upwelling and suppresses phytoplankton growth. However, such trends do not apply to this study, as some red tide-affected areas may have experienced high river discharge, allowing phytoplankton growth despite high SST. A similar pattern of positive correlation between SST and chl-a was observed in areas with high nutrient runoff. Likewise, a weak negative correlation was evident between chl-a and PAR, implying that while PAR supports photosynthesis, excessive light can trigger photoinhibition and reduce photosynthetic efficiency. Similar to chl-a,  $K_d(490)$  also shows moderate correlation with  $b_{bp}(443)$ , suggesting that more suspended particles result in higher light backscattering and attenuation. In contrast, it has weak correlations with SST and PAR, revealing that temperature and the amount of sunlight have less influence on light penetration than particles.

The weak correlation between SST and PAR has also been observed, revealing that higher sea surface temperatures do not always result in higher levels of PAR. In contrast, some

observations of increased SST and PAR levels during the dry season are accompanied by decreases in algal biomass, whereas decreases in both parameters during the wet season enhance algal growth. Meanwhile, the negligible or weak negative correlations of  $b_{bp}(443)$  with PAR and SST demonstrate its responsiveness to other parameters. Particularly, while  $b_{bp}(\lambda)$  is linked with phytoplankton carbon, the presence of non-algal particles also plays a significant role in influencing  $b_{bp}(\lambda)$ .

While there are slight variations in the findings, these only highlight the dynamic nature of marine ecosystems. These ecosystems in various water bodies can have distinct features that vary with time and space (Titaley et al., 2024). Thus, the interactions between oceanographic parameters vary depending on the region and the time being considered.

**3.1.2 Lagged Effects of Oceanographic Parameters on Red Tide:** Figure 4 shows the plot integrating 2003 red tide records with satellite-derived parameter values for bays and coastal waters. These plots were used to estimate all possible lag periods per parameter, revealing that a one (1) week lag was most frequent. On average, the five (5) parameters underwent changes a week before BFAR issued red tide advisories, affecting the timing of red tide alerts. These observations informed the creation of lag features for chl-a, SST, PAR,  $K_d(490)$ , and  $b_{bp}(443)$ . Since the parameter values were from 8-day composites, the lag features corresponded to the previous 8-day window value for each parameter, enabling the model to identify temporal patterns and historical trends and improving red tide prediction (Analytics Vidhya, 2024).

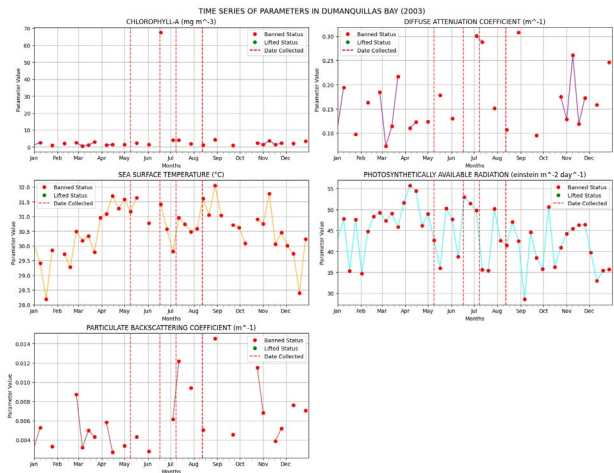


Figure 4. Parameter Time Series in Dumanquillas Bay in 2003

3.2 Time Series Analysis of Oceanographic Parameters

Most of the trends in the STL plots were fluctuating across all sites and parameters, with short-term spikes and inconsistent patterns, reflecting the realistic and dynamic environmental conditions commonly linked to red tide events. For instance, a sudden increase followed by sharp declines in chl-a suggest a bloom-bust mechanism, typical of phytoplankton dynamics. Similar trends in  $K_d(490)$  and slight deviations in  $b_{bp}(443)$  further indicate the influence of phytoplankton presence on light penetration, while non-algal particles may explain anomalies. SST trends showed a general increase over time, likely due to global warming, while PAR varied seasonally with sunlight availability.

The clustering of seasonality patterns revealed that FMAs 2, 6, 7, 9, 10, and 11 mostly followed one or two low-high cycles per year, likely influenced by monsoonal patterns. In contrast, FMAs 3, 4, and 8 showed more complex cycles, possibly due to localized anthropogenic or environmental factors. FMAs 5 and 12 showed no consistent patterns, indicating highly site-specific conditions. These seasonal trends were particularly consistent between chl-a and  $K_d(490)$ , supporting their strong correlation in understanding phytoplankton growth. Similarly,  $b_{bp}(443)$  showed recurring two low-high cycles in many FMAs, reflecting particulate matter variation likely driven by sediment resuspension during monsoon transitions, allowing phytoplankton growth. SST and PAR exhibited two high-low cycles annually, consistent with seasonal monsoons.

Overall, while seasonality is present across all sites, the pattern, timing, and frequency of these seasonal cycles vary considerably, especially in parameters heavily affected by localized environmental and anthropogenic factors, such as chl-a,  $K_d(490)$ , and  $b_{bp}(443)$ . These seasonal variations only highlight the importance of understanding context-specific environmental conditions in analyzing red tide occurrences.

3.3 Model Development

Table 2 shows that the learning rate of the model was extremely slow, preventing it from overfitting. To compensate for the low learning rate, a high number of trees ( $n_{estimators}$ ) was applied, allowing the model to pick up intricate patterns in the data over time slowly. Moderate deep trees ( $max\_depth$ ) and splits with little data ( $min\_child\_weight$ ) were also used to capture significant interactions without becoming too complicated. A considerably small value for the gamma was selected, enabling the model to be flexible and open to potential splits to improve accuracy. The subsample and colsample\_bytree values indicate how much data and features each tree used, reducing overfitting and improving generalization to newer data. For regularization, both L1 ( $reg\_alpha$ ) and L2 ( $reg\_lambda$ ) penalties were moderate, acting as barriers to control model complexity and avoid overfitting. The small  $max\_delta\_step$  allowed the model learning to be more stable by limiting the maximum allowed adjustment of the model per round of learning. Lastly, the model grew trees based on balanced depth, which can be suited for temporal and spatial oceanographic data.

Hyperparameter	Value
learning_rate	0.008727037
n_estimators	1084
max_depth	4
min_child_weight	1
gamma	8.07E-06
subsample	0.826645452
colsample_bytree	0.603955538
reg_alpha	0.194025081
reg_lambda	0.268606812
max_delta_step	2
grow_policy	depthwise

Table 2. Optimal Hyperparameters for the Model

3.4 Model Evaluation

Using the test set, the model performed moderately, as shown in Table 3, with accuracy (58%), F1-score (59%), and AUC (61%). Precision for Class 0 (lifted) and 1 (banned) was 58%, the recall was 58% and 59%, and F1 scores were also 58% and 59%, respectively. The macro and weighted averages across metrics



were 58%. It surpassed the baseline model performance produced by LazyPredict with accuracy, balanced accuracy, ROC AUC, and F1 score of 54%. These values were similar to the findings of Inokuchi et al. (2024), who applied a comparable methodology for red tide detection and achieved metrics of 58-61% using the VGG11 model. Although moderate, it performed relatively steadily and well-balanced, showing its potential as a practical solution in resource-constrained situations and as a stepping-stone for red tide early detection using satellite imagery through complementing existing methods. The non-inclusion of spatial aspects surrounding the red tide-affected sites, such as environmental factors and anthropogenic influences, may have contributed to limiting its ability to entirely capture the complex ecological interactions influencing red tide events, potentially reducing its performance.

Test Set				
Accuracy	0.58			
F1-score	0.59			
AUC	0.61			
	Precision	Recall	F1-Score	Support
Class 0 (Lifted)	0.58	0.58	0.58	493
Class 1 (Banned)	0.58	0.59	0.59	494
Macro Average	0.58	0.58	0.58	987
Weighted Average	0.58	0.58	0.58	987

Table 3. Performance Metrics for Test Set

SHAP analysis gives directional insights into the role of each oceanographic parameter in predicting red tide bans, as indicated by positive SHAP values, and lifts, as represented by negative SHAP values (Britton et al., 2024). As shown in Figure 5, SST had the strongest impact, where higher values leaned toward lifted predictions and lower values toward banned predictions, a pattern mirrored by its 8-day lag. Chl-a and its lag showed a similar, but weaker, pattern. In contrast,  $K_d(490)$  and its lagged counterpart have a smaller overall impact, with higher values leaning toward banned predictions and lower values toward lifted ones. Meanwhile, PAR and  $b_{bp}(443)$ , including their lags, showed complex and non-linear relationships with both high and low values contributing to both SHAP directions, indicating a non-linear influence on predictions.

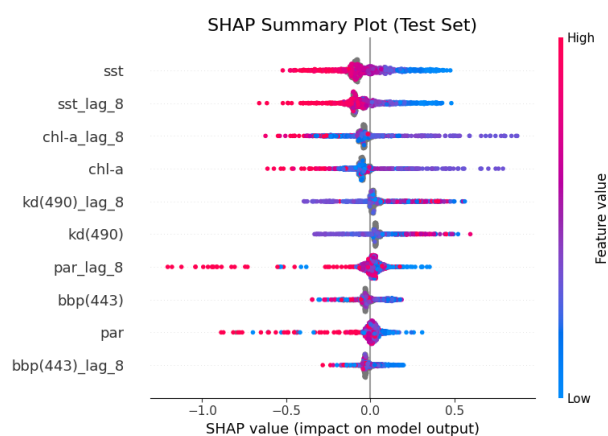


Figure 5. SHAP Summary Plot of the Test Set

These findings align with Karki et al. (2018), Hill et al. (2020), and Izadi et al. (2021), identifying SST, chl-a, and  $K_d(490)$  as key red tide detection drivers. The behavior of SST and its lag relates to the most prevalent causative organism of red tide in the Philippines, *P. bahamense*, which grows best at 28–30 °C, suggesting temperature beyond this level may limit its growth

(Smayda, 1997; Folio & Yap-Dejeto, 2022). Meanwhile, the behavior of chl-a and its lag can be explained by distinct water profiles in the Philippines, where its trends vary across areas.

### 3.5 Visualization of Prediction Results

Based on the test data, Table 4 shows overall moderate to weak model performance in FMAs 2, 5, 7, 8, and 9 across the seasons. These are attributed to similar environmental conditions across clustered red-tide affected sites, minimal differences in parameter values across seasons, missing pixel values in crucial parameters such as chl-a and  $K_d(490)$ , and training data distribution. Collectively, these reduce the ability of the model to distinguish between red tide and non-red tide events. Distinct parameter values across seasons in FMA 7 helped to yield better performance compared to other FMAs. While the results in FMAs 3 and 6 appear promising, the predictive performance of the model should be interpreted and analyzed carefully, given that each has only one site, hence, only reflecting the predictive performance at those particular sites and limiting the reliability of the findings at a broader regional scale. Poor model performance was observed in FMAs 10 and 11 due to minimal differences in parameter values across seasons and mismatched parameter values and red tide-linked conditions. Since FMA 7 had the most historical data across the seasons and the most number of sites, it has the most representative and reliable prediction results among all, with its relatively good performance metrics, followed by FMA 2 in terms of generalizability of the results.

FMA	Precision	Accuracy	Recall	F1-Score	AUC
2	0.67	0.55	0.62	0.64	0.51
3	<b>0.92</b>	<b>0.86</b>	<b>0.92</b>	<b>0.92</b>	<b>0.46</b>
5	0.38	0.68	0.74	0.50	0.70
6	<b>0.86</b>	<b>0.94</b>	<b>1.00</b>	<b>0.92</b>	<b>0.96</b>
7	0.49	0.58	0.61	0.55	0.58
8	<b>0.66</b>	<b>0.63</b>	<b>0.51</b>	<b>0.57</b>	<b>0.63</b>
9	0.39	0.63	0.73	0.51	0.66
10	0.69	0.46	0.45	0.54	0.48
11	0.48	0.50	0.50	0.49	0.50
Macro Average	0.58	0.58	0.58	0.58	0.58

Table 4. Performance Metrics across FMAs

Figures 6-9 show that correct predictions are evenly spread across FMAs in the dry season and become concentrated mainly in FMA 7 during the wet season. Meanwhile, incorrect predictions are more localized in FMA 7 during the dry season and spread out to its adjacent FMAs, particularly 2 and 10, during the wet season. This indicates that while FMA 7 had the most recorded bans for both seasons, which could help the model to predict correctly, the model still struggled to identify the red tide events in the region. Notably, the model exhibited moderate to poor predictive performance in neighboring FMAs 2, 7, and 10.

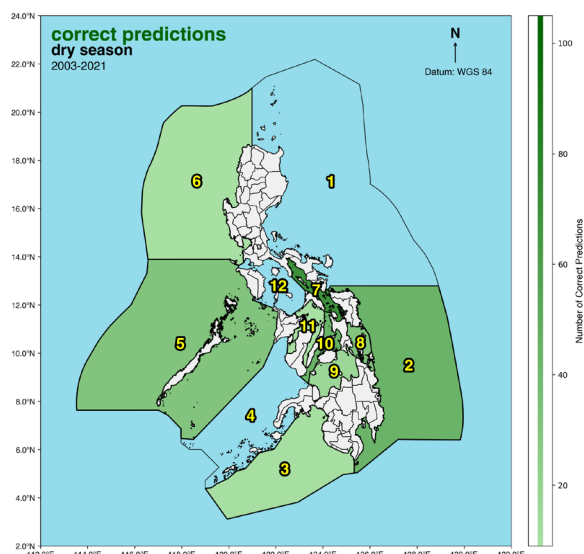


Figure 6. Spatial Distribution of Correct Predictions (Dry Season)

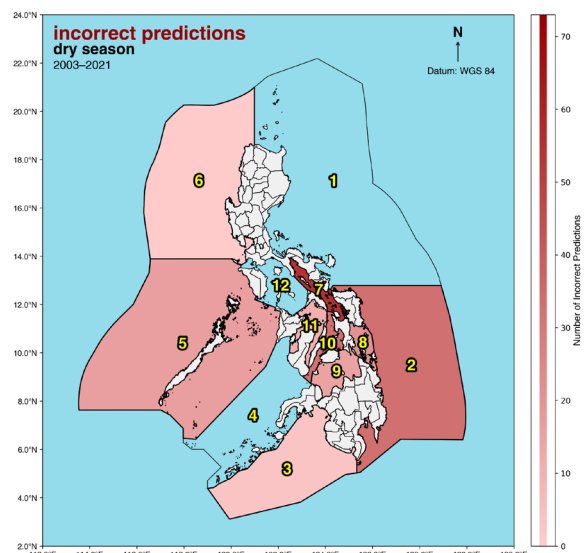


Figure 7. Spatial Distribution of Incorrect Predictions (Dry Season)

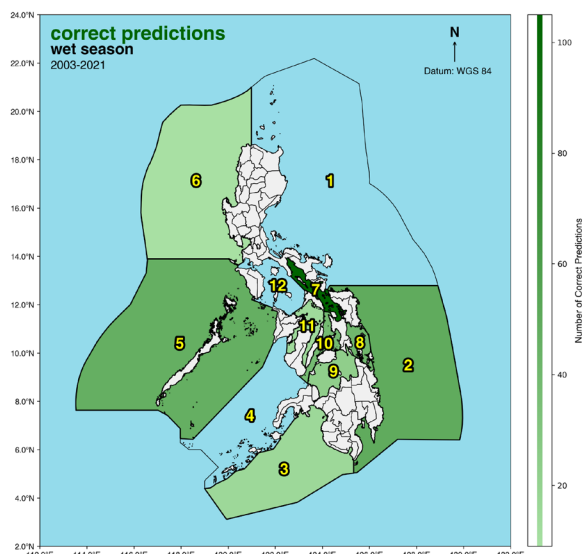


Figure 8. Spatial Distribution of Correct Predictions (Wet Season)

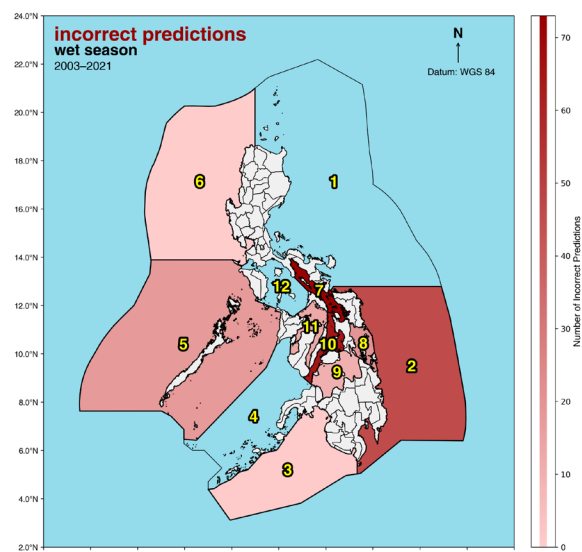


Figure 9. Spatial Distribution of Incorrect Predictions (Wet Season)

Given the spatial resolution of 4.6 km, sudden seasonal environmental variations may not easily be distinguished due to overlapping or similar patterns across adjacent FMAs. The incompleteness of parameter values may also prevent the model from recognizing red tide-linked environmental patterns. Specifically, there were instances wherein only one or two out of the five parameters had valid values. Altogether, these affect the predictive accuracy of the model across FMAs and seasons.

#### 4. Conclusions and Recommendations

This study developed a predictive model for red tide events in the Philippines using MODIS-derived oceanographic parameters and XGBoost. The analysis revealed a strong positive correlation between chl-a and  $K_d(490)$ , which indicates that increased phytoplankton biomass inhibits light penetration in the water column, making them strong predictors of red tide. The moderate correlation between chl-a and  $b_{bp}(443)$  still supports this notion, indicating the presence of both algal and non-algal particles. Meanwhile, weak correlations of SST and PAR can be attributed to excessive heat and light, causing delayed phytoplankton growth and photoinhibition. Based on STL, chl-a and  $K_d(490)$  shared similar and consistent seasonal and trend behavior across red tide-affected areas. All parameters followed seasonal patterns, with most showing fluctuating trends over time, which reflects the dynamic nature of red tide and the influence of localized factors like nutrient runoff and mariculture. These suggest that red tide detection should account not only for parameter interactions but also for their temporal trends and the localized environmental conditions. The model achieved moderate and stable performance as reflected in 58-59% accuracy and F1-score, and 61% AUC with SST, chl-a, and  $K_d(490)$  as key red tide predictors. While this is not a replacement for existing methods, it serves as a support tool and foundation for early detection using MODIS-derived data. Model performance across different FMAs and seasons varied due to factors such as minor variation in parameter values across adjacent FMAs and seasons, missing pixel values of crucial parameters, mismatched parameter values, red tide-linked conditions, incomplete and unevenly distributed training data across seasons. Among all, the model produced the most reliable

and representative results in FMA 7, and the poorest in FMAs 10 and 11.

Based on these findings, it is recommended to use higher-resolution satellite imagery to better capture the spatial variability in the oceanographic parameters and refine the prediction results. Incorporating additional satellite-derived oceanographic and external environmental parameters, and including cell density or toxin level data for regression-based modeling when available, may also improve the results.

## References

- Azanza, R. V., Yñiguez, A. T., Onda, D. F., Benico, G. A., Lim, P. T., Leaw, C. P., & Iwataki, M., 2024. Expansion of toxic algal blooms in coastal and marine areas in the Philippines and Malaysia: Is it climate change related? *Sustainability*, 16(8), 3304. doi.org/10.3390/su16083304.
- Britton, A., Graham, G., & Woloszyn, M., 2024. Exploring Relationships Between Drought Indices and Ecological Drought Impacts Using Machine Learning and Explainable AI. *Journal of Applied and Service Climatology*, 2024(005), 1–12. doi.org/10.46275/joasc.2024.09.001.
- Ching, P. K., Ventura, R. J., de los Reyes, V. C., Sucaldito, M. N., & Tayag, E., 2015. Lethal paralytic shellfish poisoning from consumption of green mussel broth, Western Samar, Philippines, August 2013. *Western Pacific Surveillance and Response Journal*, 6(2), 22–26. doi.org/10.5365/wpsar.2015.6.1.004.
- Delisi, J., 2024. How to Prepare Data for Use in Machine Learning Models. PhData. phdata.io/blog/how-to-prepare-data-for-use-in-machine-learning-models (11 August 2025).
- Donatelli, R. E., Ji Ae Park, Mathews, S. M., & Lee, S.-J., 2022. Time series analysis. *American Journal of Orthodontics and Dentofacial Orthopedics*, 161(4), 605–608. doi.org/10.1016/j.ajodo.2021.07.013.
- Folio, F. M., & Yap-Dejeto, L., 2022. Phytoplankton Composition during a Period of the Red Tide Bans in 2017 in Irong-Irong Bay, Western Samar, Philippines. *Philippine Journal of Science*, 151(S1). doi.org/10.56899/151.s1.16.
- Gordan, M.-I., Popescu, C. A., Călina, J., Adamov, T. C., Mănescu, C. M., & Iancu, T., 2024. Spatial Analysis of Seasonal and Trend Patterns in Romanian Agritourism Arrivals Using Seasonal-Trend Decomposition Using LOESS. *Agriculture*, 14(2), 229. doi.org/10.3390/agriculture14020229.
- Hill, P. R., Kumar, A., Temimi, M., & Bull, D. R., 2020. HABNet: Machine Learning, Remote Sensing-Based Detection of Harmful Algal Blooms. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 13, 3229–3239. doi.org/10.1109/jstars.2020.3001445.
- Inokuchi, Y., Kobayashi, K., Guillerault, J., Henmi, Y., Gonzalez, P. H., Aritsugi, M., & Mendonca, I., 2024. Development of a Red Tide Early Detection System Using Satellite Images. *IEEE*. doi.org/10.1109/gecost60902.2024.10474956.
- Izadi, M., Sultan, M., Kadiri, R. E., Ghannadi, A., & Abdelmohsen, K., 2021. A Remote Sensing and Machine Learning-Based Approach to Forecast the Onset of Harmful Algal Bloom. *Remote Sensing*, 13(19), 3863. doi.org/10.3390/rs13193863.
- Joshi, N., Park, J., Zhao, K., Londo, A., & Sami Khanal., 2024. Monitoring Harmful Algal Blooms and Water Quality Using Sentinel-3 OLCI Satellite Imagery with Machine Learning. *Remote Sensing*, 16(13), 2444–2444. doi.org/10.3390/rs16132444.
- Karki, S., Sultan, M., Elkadiri, R., & Elbayoumi, T., 2018. Mapping and Forecasting Onsets of Harmful Algal Blooms Using MODIS Data over Coastal Waters Surrounding Charlotte County, Florida. *Remote Sensing*, 10(10), 1656. doi.org/10.3390/rs10101656.
- Khan, R. M., Salehi, B., Mahdianpari, M., Mohammadimanesh, F., Mountrakis, G., & Quackenbush, L. J., 2021. A Meta-Analysis on Harmful Algal Bloom (HAB) Detection and Monitoring: A Remote Sensing Perspective. *Remote Sensing*, 13(21), 4347. doi.org/10.3390/rs13214347.
- Manzar Abbas, M., Melesse, A. M., Scinto, L. J., & Rehage, J. S., 2019. Satellite Estimation of Chlorophyll-a Using Moderate Resolution Imaging Spectroradiometer (MODIS) Sensor in Shallow Coastal Water Bodies: Validation and Improvement. *Water*, 11(8), 1621. doi.org/10.3390/w11081621.
- Pal, M., Yesankar, P. J., Dwivedi, A., & Qureshi, A., 2020. Biotic control of harmful algal blooms (HABs): A brief review. *Journal of Environmental Management*, 268(110687), 110687. doi.org/10.1016/j.jenvman.2020.110687.
- Smayda, T.J., 1997. Harmful algal blooms: Their ecophysiology and general relevance to phytoplankton blooms in the sea. *Limnology and Oceanography*, 42, 1137–1153. doi.org/10.4319/lo.1997.42.5\_part\_2.1137.
- Titaley, J., Lumingas, L.J.L., Pinontoan, B., Nanlohy, P., & Manu, L., 2024. Spatial and Temporal Analysis of Oceanographic Parameters and Their Relationship with Upwelling Phenomena in Seram Sea and Buru Island Waters. *Egyptian Journal of Aquatic Biology and Fisheries*, 28(5), 379–397. doi.org/10.21608/ejabf.2024.378880.
- Wexler, P., 2014. Encyclopedia of Toxicology | ScienceDirect. sciencedirect.com/referencework/9780123864550/encyclopedia-of-toxicology (11 August 2025).
- Yñiguez, A. T., Lim, P. T., Leaw, C. P., Jipanin, S. J., Iwataki, M., Benico, G., & Azanza, R. V., 2020. Over 30 years of HABs in the Philippines and Malaysia: What have we learned? *Harmful Algae*, 102, 101776. doi.org/10.1016/j.hal.2020.101776.