

Honours Thesis

Modelling Coccolithophore data in the North Atlantic Ocean

Seoyeon (Cali) Park

B00768397

Fall 2021

# Contents

<b>1</b>	<b>Abstract</b>	<b>3</b>
<b>2</b>	<b>Acknowledgements</b>	<b>3</b>
<b>3</b>	<b>Introduction</b>	<b>4</b>
<b>4</b>	<b>Methods</b>	<b>5</b>
4.1	Data collection and manipulation . . . . .	5
4.2	Brief introduction of generalized additive models (GAM) and hierarchical generalized linear models (HGLM) . . . . .	7
4.3	Hierarchical Generalized Additive Models (HGAM) . . . . .	8
4.4	Least Absolute Shrinkage and Selection Operator (LASSO) . . . . .	9
4.5	Software . . . . .	10
4.6	Training and Test data . . . . .	10
4.7	Model Comparison Methods . . . . .	11
<b>5</b>	<b>Results</b>	<b>11</b>
5.1	HGAMS . . . . .	11
5.1.1	R-squared adjusted values . . . . .	11
5.1.2	AIC values . . . . .	12
5.1.3	RMSE values . . . . .	13
5.2	LASSO regression . . . . .	14
5.2.1	Coefficients of terms . . . . .	14
5.2.2	RMSE values . . . . .	16
<b>6</b>	<b>Discussion and conclusions</b>	<b>18</b>
<b>7</b>	<b>Appendix</b>	<b>21</b>

# 1 Abstract

This thesis builds species models of coccolithophores using hierarchical generalized additive models (HGAM) and regression analysis using Least Absolute Shrinkage and Selection Operator (LASSO). Analysis will include many common model comparison techniques such as comparisons of adjusted  $R^2$  and Akaiake Information Criterion scores of the training data, and Root Mean Squared Errors of both training and testing data, to examine which models can accurately predict species abundance. We start by extracting a subset of the Coccobase data repository and apply a series of models and analysis on corresponding results. Three of the best HGAMs were chosen based on comparison techniques mentioned above. The best LASSO model assessed was for the *Emiliania huxleyi* data. These components will be interpreted and discussed in depth.

# 2 Acknowledgements

I would like to express my gratefulness to all individuals who have contributed to aid in writing this honours thesis. I would like to acknowledge Dr. Andrew Irwin, whom without his support and patience would not have been possible. Meeting every week for discussions on room to improve and how to approach issues that have arisen during this process have been a huge support. Without Dr. Irwin, I would not have been able to gain knowledge and experience with applied statistical research with oceanography data.

I would like to thank the lab members, especially Joe, whom have initially taught me the background of LASSO regression, and Mohammad, always answering questions about statistics, both very kind and knowledgeable in their research. Another thanks to Crispin, who also suggested different approaches to summarizing the HGAM smooths, when I was stuck with a problem.

Special thanks to the professors that have taught me statistics throughout my undergraduate degree. I would like to thank each and every person that have taught me the world of statistics. Having a solid foundation from their shared knowledge have allowed me to want to pursue a further Master's degree in Statistics.

### 3 Introduction

Coccolithophores are a major type of calcifying phytoplankton commonly found across all oceans that occupy distinct niches from tropical to sub-polar seas. As climatic conditions change across the oceans from various reasons including global warming, it is important to understand how species abundance will change, and to develop species distribution models to predict how their future biogeography may change as a consequence of anticipated changes in climate. Many previous studies, using controlled laboratory experiments, show that variables such as light, nutrients, and temperature play an important role in determining interspecific competition and species abundance. Analysis of species abundance and its relation to environmental and physical variables using real data collected across the world’s oceans is the next necessary step in accurately predicting species abundances and finding similarities and/or differences in the environmental conditions where coccolithophore species are observed, and is the main objective of this work.

Coccobase is a compilation of previously reported data that include both environmental and physical observations of coccolithophore species across the globe as well as their respective abundances [6]. Coccobase therefore provides an excellent opportunity to model coccolithophore abundances in the North Atlantic using real ocean data, which may in the future be extrapolated to the rest of the Earth’s oceans.

Species models for *Emiliana huxleyi*, *Coccolithus pelagicus*, *Syracosphaera pulchra* are compared to models for the broad category of all coccolithophore observations, and the genera *Coccolithus* and *Syracosphaera*. These species were chosen based on the number of observations available. *E. huxleyi*, *C. pelagicus*, *S. pulchra* were the species with the highest number of observations after subsetting the data (Table 11).

The purpose of this research was to determine which environmental variables can be used to accurately predict species abundance, and to find similarities and/or differences in the environmental conditions where the coccolithophore species are observed, using appropriate modelling methods. The data was modelled using hierarchical generalized additive models (HGAMs) and GLMs with LASSO variable selection and regularization [3, 9]. These models were used to assess the differences in niche and to compare the predictive ability of the different modelling techniques. The best model for predictive accuracy and fit of data were determined based on common model comparison methods, using values of adjusted  $R^2$  (R-sq adj), Akaike information criterion scores (AIC), and root mean squared errors (RMSE) [1].

## 4 Methods

### 4.1 Data collection and manipulation

The data used for this analysis comes from the Coccobase repository, available on Github [4]. Coccobase is built upon previously made compilations of data and includes many of the variables required for analysis. Coccobase has over 213,000 individual observations, carefully filtered by applying quality control measures to the pre-existing data gathered from the literature and earlier databases. There are 190 different coccolithophore taxa reported in Coccobase, but for this report, a subset of the data will be used [6].

The data included all observations of coccolithophores in the North Atlantic. (Table 11) The models include observations of coccolithophores, of the genera *Coccolithus* and *Syracosphaera*, as well as species specific observations of *Emiliana huxleyi*, *Coccolithus pelagicus*, and *Syracosphaera pulchra*. All the species-specific observations were subsetting based on the highest number of observations. The following table describes the number of observations in data corresponding to each species. ea

Table 1: Species and corresponding number of observations

Family/Species name	Number of observations
<i>Coccolithophore</i>	5823
<i>Coccolithus</i> species	803
<i>Syracosphaera</i> species	790
<i>Emiliana huxleyi</i>	1367
<i>Coccolithus pelagicus</i>	572
<i>Syracosphaera pulchra</i>	240

The variables used to predict species abundance were: phosphate, Si\*, N\*, salinity, temperature, and a discrete classifier for each species. Phosphate, is an essential element for nucleic acid, RNA/DNA synthesis of coccolithophores. [7] The higher the phosphate concentration, the more beneficial it would be for the species. Si\* is the difference between silicate and nitrate concentration, and N\* is the difference between nitrate and 16 times the phosphate concentration. Both differences were chosen because they show how species abundance

and species competition are correlated to the availability of environmental variables. Temperature is known to affect the abundance of coccolithophore species, hence it is sufficient to include it as one of the predictor variables [10]. Finally, salinity is a predictor because it influences strain specific growth. These variables will require an addition of the tensor product between longitude and latitude in the North Atlantic Ocean. Observations for species abundance was log transformed to follow an approximate normal distribution, since the histogram of the original scale appeared to be positively skewed.

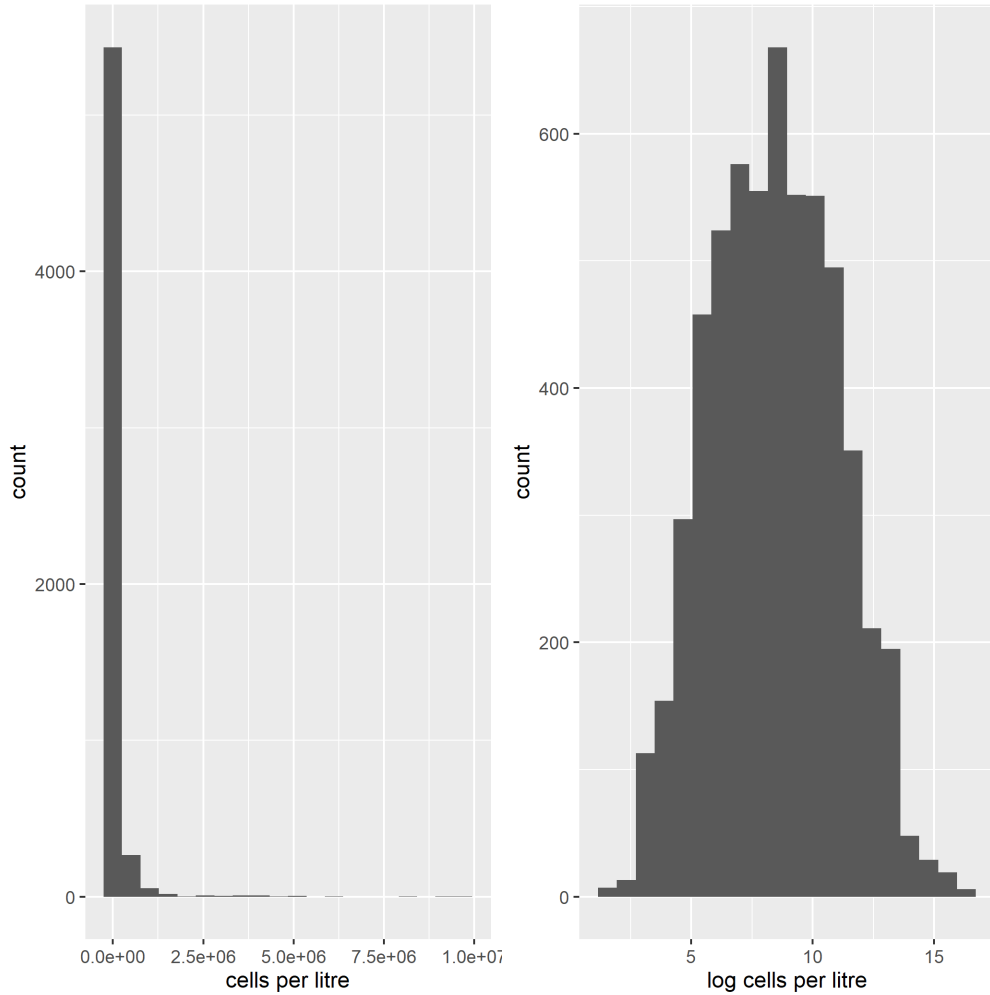


Figure 1: Histogram of species abundance counts. Left: Original scale. Right: Log transformed.

## 4.2 Brief introduction of generalized additive models (GAM) and hierarchical generalized linear models (HGLM)

Two modelling techniques, Generalized Additive Models (GAM) and Hierarchical Generalized Linear Models (HGLM), have been commonly used to model highly variable data. Certain aspects of these models are useful to construct the Hierarchical Generalized Additive Models, which was one of the two modelling techniques used in this study.

GAMs are very flexible models, where the only assumptions are that the functions are smooth (via splines) and that the response is an additive function of the smooths of each predictor. They have the following structure:

$$E[Y] = g^{-1}(\beta_0 + \sum_{i=1}^I f_i(x_i)) \quad (1)$$

where  $E(Y)$  is the expected value of the response  $Y$  with the corresponding link function  $g$ , and  $f_i(x_i)$  is the sum of all predictors smoothed by the function  $f_i$  with intercept  $\beta_0$ . The response is related upon the linear combination of smoothed  $(s_i(x_i), i = 1, \dots, I)$  predictor variables, which usually allows for better fit of the data [3].

One of the advantages of using GAMs is that they are highly suitable for fitting non-linear data. As mentioned above, there are many ways to smooth the predictor variables via splines. Types of splines include Thin Plate Regression Splines (TPRS) and Cyclic Cubic Regression Splines (CRS), which should be chosen carefully based on the nature of the predictor variables.

Table 2: Description of "smoothers" used in analysis.

Smoothers	Description
Thin Plate Regression Splines (TPRS)	Used to smooth covariates in any number of dimensions
Cyclic Cubic Regression Splines (CRS)	Useful for fitting covariates that undergo cyclic effects (seasonal)
Random Effects (RE)	Smooths covariates without any assumptions on the structure.

In the analysis, phosphate, N\*, Si\*, salinity, temperature, and latitude used this smoother TPRS (although the temperature variable may suffice using the CRS smoother), where as longitude and scientific name used

the smoothers CRS and RE, respectively.

Another commonly used modelling technique, are HGLMs (or also called generalized linear mixed models) which are a flexible extension of generalized linear models (GLMs) [5]. They meet the assumptions of a generalized linear model with an additional component that allows the model to have hierarchical (structural or multi-level) relationships between the predictors, and as well have an overall relationship between the predictors and the response. This hierarchical structure leads to building more realistic models and reduces overfitting [2]. The assumptions of these modelling techniques will be beneficial to keep in mind for the development of the HGAM, explained in the next section.

### 4.3 Hierarchical Generalized Additive Models (HGAM)

Components of HGLM and GAM can be combined to make hierarchical generalized additive models (HGAM). HGAMs are constructed to model highly variable data, to understand the smooth hierarchical functional relationships between predictor groups and the response. The modelling efforts of this report closely followed the methods outlined in Pederson et al [9]. For this paper, only models G, GS, and GI are performed for analysis of coccolithophores in the North Atlantic Ocean. Model G are fitted on datasets including all coccolithophore observations, and datasets including the *Coccolithus* and *Syracosphaera* species. Model GI and GS are fitted on species specific data subsets. All models are fitted using the restricted (residual) maximum likelihood method.

Description of each model is summarized in the following table:

Table 3: Types of Models in Analysis

Type of Model	Description
Model G	A single global smoother for all observations
Model GS	A single global smoother plus group-level smoothers that have the same wigginess
Model GI	A single global smoother plus group-level smoothers that have different wigginess

Each model has the following structure coded into R:



- Model G: `gam(cells_per_litre = s(p_an, bs = "tp") + s(si_star, bs = "tp") + s(ni_star, bs = "tp") + s(s_an, bs = "tp") + s(t_an, bs = "tp") + s(scientific_name, bs = "re") + te(longitude, latitude, bs = c("cc","tp"))), data)`
- Model GS (phosphate): `gam(cells_per_litre = s(p_an, bs = "tp") + s(si_star, bs = "tp") + s(p_an, bs = "tp") + s(s_an, bs = "tp") + s(t_an, bs = "tp") + te(longitude, latitude, bs = c("cc", "tp"))) + s(p_an, scientific_name, bs = "fs"), data)`
- Model GI (phosphate): `gam(cells_per_litre = s(p_an, bs = "tp") + s(si_star, bs = "tp") + s(p_an, bs = "tp", m = 1) + s(s_an, bs = "tp") + s(t_an, bs = "tp") + te(longitude, latitude, bs = c("cc", "tp"))) + s(p_an, scientific_name, bs = "fs"), data)`

#### 4.4 Least Absolute Shrinkage and Selection Operator (LASSO)

The Least Absolute Shrinkage and Selection Operator (LASSO) can be incorporated into linear models to shrink some estimated parameters to 0. This allows for automatic variable selection and regularization for building simpler and more interpretable models.

Suppose the data is composed of  $N$  total observations. Under the assumption that each observation is independent (the same assumptions as ordinary least squares, OLS), the outcome of the data is arranged into an  $n \times 1$  vector:

$$y_i = \begin{bmatrix} y_1 & y_2 & \dots & y_n \end{bmatrix}^\top$$

where  $i = 1, \dots, n$ .

Each predictor,  $x_i$ , is arranged into a  $n \times p$  data matrix of observations:

$$x_{np} = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \dots & \dots & \dots & \dots \\ x_{n1} & x_{n2} & \dots & x_{np} \end{bmatrix}$$

where  $n$  is number of predictor variables in the data, and  $p$  is the total number of observations corresponding to each predictor variable.

In addition, suppose that there is an  $n \times 1$  vector of parameters for each predictor  $\beta_n$ :

$$\hat{\beta} = \begin{bmatrix} \beta_1 & \beta_2 & \dots & \beta_n \end{bmatrix}^\top$$

The standard linear model is:

$$y = X\beta + \epsilon$$

LASSO regression then becomes a modified least squares problem:

$$\hat{\beta}_\lambda = \underset{\beta}{\operatorname{argmin}} \left\| Y - X\beta \right\|_2^2 + \lambda \left\| \beta \right\|_1$$

where  $\left\| \beta \right\|_1 = \sum_{i=1} |\beta_i|$ . The addition of the  $\lambda \left\| \beta \right\|_1$  term adds regularization to the least squares problem, by adding a penalty to the minimization of  $\left\| Y - X\beta \right\|_2^2$ .

In addition, we see that for any  $\lambda > 0$ , there exists an  $s_\lambda = \left\| \beta_\lambda \right\|_1$  where  $\hat{\beta}_\lambda = \underset{\beta}{\operatorname{argmin}} \left\| Y - X\beta \right\|_2^2$  such that  $\|\beta_i\| \leq s_\lambda$  [11].

## 4.5 Software

All data is manipulated and modelled using pre-existing R packages. Refer to the appendix for a full list of packages used. The coding was performed with R/RStudio. The full R project can be accessed via Github at: <https://github.com/calipark1213/Honours-Coccobase> [8].

## 4.6 Training and Test data

Each of the species subsets were further partitioned into training and testing data. To allow both sets to have a good representation of data and to avoid overfitting, total data were split randomly into 60% training and 40% test data. The number of observations outlined in table 1 are further split into the following table:

Table 4: Training and Testing set

Genus/Species name	Train	Test
<i>Coccolithophore</i>	3493	2329
<i>Coccolithus</i> species	482	321
<i>Syracosphaera</i> species	474	316
<i>Emiliana huxleyi</i>	820	547
<i>Coccolithus pelagicus</i>	343	229
<i>Syracosphaera pulchra</i>	144	96

## 4.7 Model Comparison Methods

Common model comparison techniques was used to analyze the differences between HGAM and LASSO regression models. The R-squared adjusted values and AIC values of the training data will be compared. R-squared adjusted is a measure of how much the variation of the full data is measured by the model, whereas AIC values evaluates how well a model fits the data it was generated from. Furthermore, the RMSE values of both training and testing data will be compared, since it measures the standard deviation of the residuals. For LASSO, models will be compared using variable selection and RMSE values. RMSE values were calculated on a log transformed scale, so the model values should be interpreted by taking the exponent of the RMSE value, and be interpreted as a factor out on average from the true values in the original scale. By taking the value on a log scale, the values are less sensitive to large outliers.

# 5 Results

## 5.1 HGAMS

### 5.1.1 R-squared adjusted values

As outlined in table 5, for models G and GS, the highest R-squared adjusted value corresponded to the dataset that included all of the *Coccolithus* genus. The lowest value of R-squared adjusted for model G

corresponded to the species specific dataset for *Syracosphaera pulchra*, with a value of 0.497 (or deviating by a factor of  $e^{0.497} = 1.6438$  on the original scale). The second lowest R-squared adjusted value was for the dataset that included all of the species of coccolithophores in the North Atlantic.

For models GS and GI, the group-level smoothers with same/different wiggleness, respectively, did not differ by a significant amount, for species *Emiliana huxleyi* and *Coccolithus pelagicus*. For model GI, the highest R-squared value was for the model with the individual penalty on the salinity variable for the *Coccolithus pelagicus* species.

Table 5: Adjusted R-squared values of all models (Training data)

Model G					
<i>Coccolithophore</i>	0.591				
<i>Coccolithus</i> species	0.718				
<i>Syracosphaera</i> species	0.641				
<i>Emiliana huxleyi</i>	0.643				
<i>Coccolithus pelagicus</i>	0.674				
<i>Syracosphaera pulchra</i>	0.497				
Model GS	Phosphate	Si*	N*	Salinity	Temperature
<i>Emiliana huxleyi</i>	0.633	0.633	0.633	0.633	0.633
<i>Coccolithus pelagicus</i>	0.639	0.639	0.639	0.639	0.639
<i>Syracosphaera pulchra</i>	0.492	0.492	0.492	0.492	0.492
Model GI	Phosphate	Si*	N*	Salinity	Temperature
<i>Emiliana huxleyi</i>	0.634	0.633	0.639	0.633	0.633
<i>Coccolithus pelagicus</i>	0.639	0.639	0.637	0.65	0.641
<i>Syracosphaera pulchra</i>	0.492	0.492	0.505	0.492	0.496

### 5.1.2 AIC values

The best models out of all the models computed for model G, GS, and GI, were the models that included the training data corresponding to the species *Emiliana huxleyi*. As outlined in table 6, the lowest AIC scores corresponded to the *Emiliana huxleyi* models. For model G, this value was -591.8112, and for model GS, the lowest AIC value corresponded to the group level smooth of salinity with *Emiliana huxleyi* as the species.

Lastly, for model GI, the lowest AIC value corresponded to the group level smooth of of N\* with *Emiliana huxleyi* as the species.

Table 6: AIC values of all models (Training data)

Model G					
<i>Coccolithophore</i>	-425.736				
<i>Coccolithus</i> species	-175.147				
<i>Syracosphaera</i> species	-351.161				
<i>Emiliana huxleyi</i>	-591.811				
<i>Coccolithus pelagicus</i>	-90.676				
<i>Syracosphaera pulchra</i>	-81.999				
Model GS	Phosphate	Si*	N*	Salinity	Temperature
<i>Emiliana huxleyi</i>	-570.196	-569.632	-571.776	-571.850	-569.631
<i>Coccolithus pelagicus</i>	-62.0486	-62.049	-62.049	-62.049	-62.049
<i>Syracosphaera pulchra</i>	-77.033	-77.032	-77.034	-77.032	-77.033
Model GI	Phosphate	Si*	N*	Salinity	Temperature
<i>Emiliana huxleyi</i>	-573.891	-569.844	-587.358	-571.846	-570.704
<i>Coccolithus pelagicus</i>	-62.049	-62.049	-59.2310	-48.448	-65.229
<i>Syracosphaera pulchra</i>	-75.423	-77.033	-77.759	-77.033	-77.173

### 5.1.3 RMSE values

Finally, to measure how well the fitted training models will predict in general, the RMSE of the predictions based on the test data were calculated. The RMSE values shared the same conclusion as the AIC values.

The best models out of all the models computed for model G, GS, and GI, were the models that included the testing data corresponding to the species *Emiliana huxleyi*. As outlined in table 7, the lowest RMSE scores for the *Emiliana huxleyi* models corresponded to the smallest AIC value. For model G, this value was 0.18261, and for model GS, the lowest AIC value corresponded to the group level smooth of phosphate with *Emiliana huxleyi* as the species. Same as the conclusion for the AIC values, for model GI, the lowest RMSE value corresponded to the group level smooth of of N\* with *Emiliana huxleyi* as the species.

Table 7: RMSE values of all models ([Training data](#), Test data)

Model G					
<i>Coccolithophore</i>	0.220	0.227			
<i>Coccolithus</i> species	0.183	0.246			
<i>Syracosphaera</i> species	0.149	0.296			
<i>Emiliana huxleyi</i>	0.159	0.183			
<i>Coccolithus pelagicus</i>	0.186	0.226			
<i>Syracosphaera pulchra</i>	0.164	0.228			
Model GS	Phosphate	Si*	N*	Salinity	Temperature
<i>Emiliana huxleyi</i>	0.163 0.182	0.163 0.182	0.163 0.182	0.163 0.182	0.163 0.182
<i>Coccolithus pelagicus</i>	0.201 0.216	0.201 0.216	0.201 0.216	0.201 0.216	0.201 0.216
<i>Syracosphaera pulchra</i>	0.165 0.221	0.165 0.221	0.165 0.226	0.165 0.221	0.165 0.221
Model GI	Phosphate	Si*	N*	Salinity	Temperature
<i>Emiliana huxleyi</i>	0.162 0.181	0.163 0.182	0.161 0.180	0.163 0.182	0.162 0.181
<i>Coccolithus pelagicus</i>	0.201 0.216	0.201 0.214	0.202 0.214	0.197 0.214	0.201 0.217
<i>Syracosphaera pulchra</i>	0.164 0.197	0.165 0.221	0.161 0.164	0.165 0.221	0.164 0.187

The top 3 models based on the R-squared adjusted, AIC, and RMSE value were: Model G of *Coccolithus* species, Model GS with the group-level smooth of salinity for *Emiliana huxleyi*, and model GI with the group-level smooth of N\* for *Emiliana huxleyi*.

## 5.2 LASSO regression

The same subsetted data were fitted using a different modelling technique, using LASSO. The hopes were that the best models based on HGAMs would also be the best choice using LASSO regression. Using the glmnet package from R, the LASSO models were fitted using the best lambda value calculated by the lambda.min function. From there, the coefficients of the model estimates were extracted. The top 3 LASSO models will also be chosen based on variable selection and RMSE values. The results are outlined below.

### 5.2.1 Coefficients of terms

LASSO models using *Coccolithus* species, *Emiliana huxleyi*, *Coccolithus pelagicus*, and the full data shrunk one of the environmental variables to 0. This is seen in table 8, where the models for *Coccolithus* and full data have shrunk the salinity term to 0, and the models for *Emiliana huxleyi* and *Coccolithus pelagicus* data

have shrunk the scientific name term to 0. The model for *Syracosphaera pulchra* have shrunk 3 terms to 0 with the coefficients being closest to 0, but with sample size taken into account this model should not be included in the final selection. For each LASSO model, the lambda.min function was used to extract the best lambda value. The following plot shows  $\text{Log}(\lambda)$  against the Mean-Squared Error component. Note that (1) - (6) is the plot for *Coccolithophore*, *Coccolithus* species, *Syracosphaera* species, *Emiliana huxleyi*, *Coccolithus pelagicus*, and *Syracosphaera pulchra*, respectively.

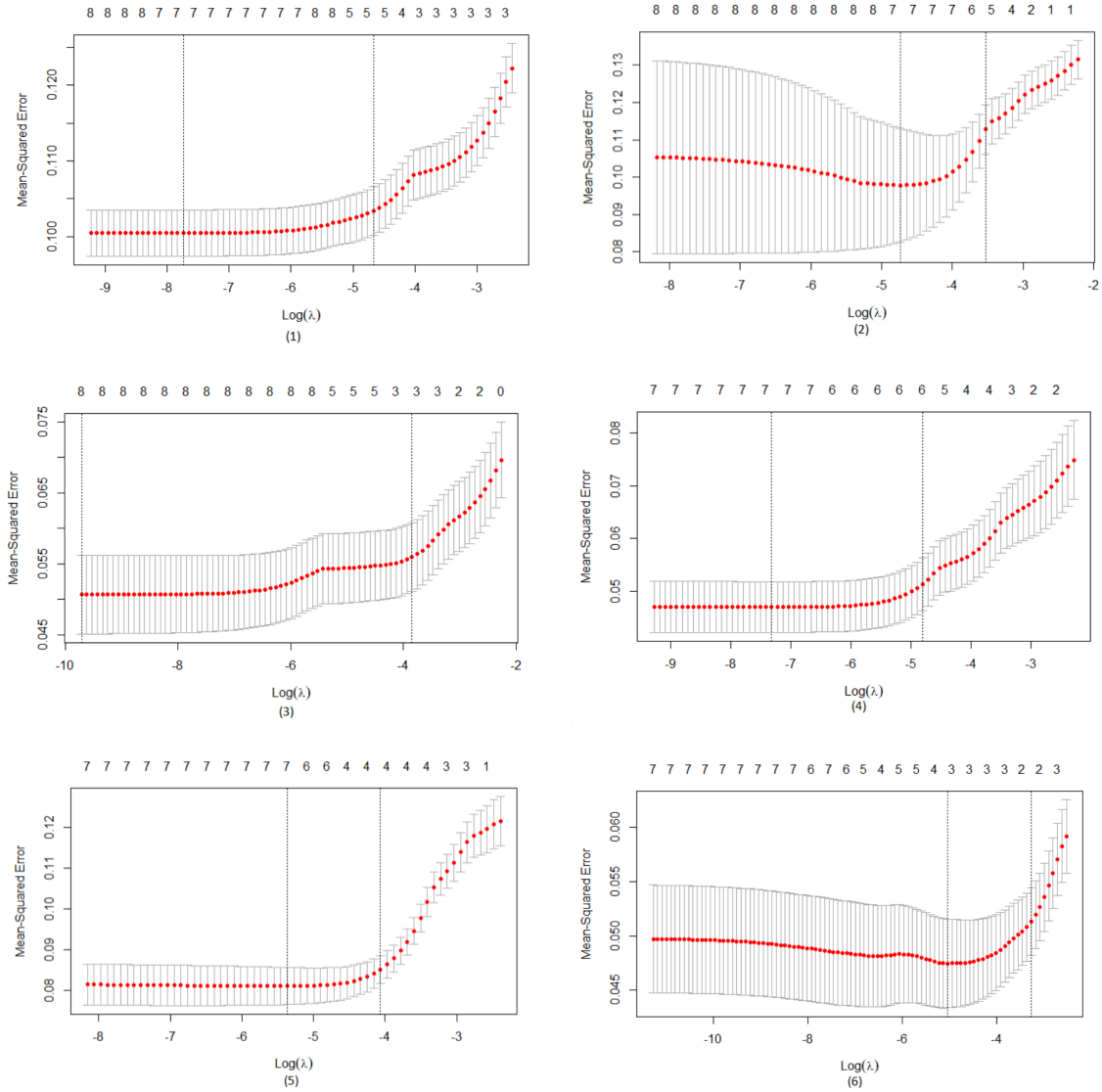


Figure 2: Plot of  $\text{Log}(\lambda)$  against Mean-Squared Error corresponding to each LASSO model.

Table 8: Table of coefficients via LASSO regression (Training data)

	<i>Coccolithophore</i>	<i>Coccolithus</i> species	<i>Syracosphaera</i> species	<i>Emiliana huxleyi</i>	<i>Coccolithus pelagicus</i>	<i>Syracosphaera pulchra</i>
Intercept	0.486	-2.493	-26.848	-7.833	-1.722	-1.833E1
Phosphate	-0.487	-0.276	-1.606	-1.231	-0.097	.
Si*	-0.087	-0.073	-0.172	-0.108	-0.065	-8.699E-2
N*	-0.039	0.012	-0.159	-0.088	0.014	-1.609E-1
Salinity	.	.	0.781	0.268	-0.141	5.572E-1
Temperature	0.039	0.069	-0.024	-0.005	0.097	2.395E-5
Scientific name	-0.003	0.010	-0.001	.	.	.
Longitude	0.004	-0.013	-0.019	0.009	-0.012	-.7.292E-3
Latitude	0.022	0.055	0.021	0.018	0.065	.

Table 9: Minimum lambda value corresponding to each model fitted with training data

Family/Species name	Minimum lambda value
<i>Coccolithophore</i>	3.9716E-4
<i>Coccolithus</i> species	9.616E-3
<i>Syracosphaera</i> species	1.2768E-4
<i>Emiliana huxleyi</i>	9.386E-5
<i>Coccolithus pelagicus</i>	7.276E-4
<i>Syracosphaera pulchra</i>	1.341E-3

Since the models were fitted using the `cv.glmnet` function, the minimum lambda (or best lambda) value were calculated for each of the regression models. The above table shows the values corresponding to each species data.

### 5.2.2 RMSE values

The following table summarizes the RMSE values of each test data fitted by LASSO regression:



Table 10: RMSE values of all LASSO models ([Training data](#), Test data)

Family/Species name	RMSE
<i>Coccolithophore</i>	<a href="#">0.316</a> 0.313
<i>Coccolithus</i> species	<a href="#">0.287</a> 0.308
<i>Syracosphaera</i> species	<a href="#">0.218</a> 0.245
<i>Emiliana huxleyi</i>	<a href="#">0.213</a> 0.223
<i>Coccolithus pelagicus</i>	<a href="#">0.275</a> 0.288
<i>Syracosphaera pulchra</i>	<a href="#">0.207</a> 0.205

From this table, we can see that the lowest RMSE value corresponds to the LASSO regression model fitted using the *Syracosphaera pulchra* data, although one should note that due to the small number of observations and the LASSO regression model shrinking most of the coefficients to zero, the next lowest RMSE value should be chosen. Therefore, the best model via LASSO regression is the model fitted onto the *Emiliana huxleyi* data.

Combining the conclusions from HGAMs and LASSO regression, the best models based on the comparison techniques outlined above were all using the *Emiliana huxleyi* data. Actual vs predicted plots of the best models are shown below:

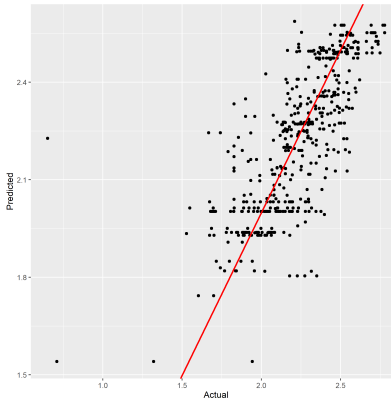


Figure 3: Actual VS Predicted model of *Emiliana huxleyi* data. Model GS with group-level smooth on salinity.

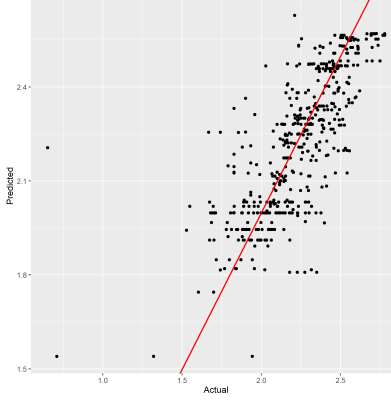


Figure 4: Actual VS Predicted model of *Emilia huxleyi* data. Model GI with group-level smooth on  $N^*$ .

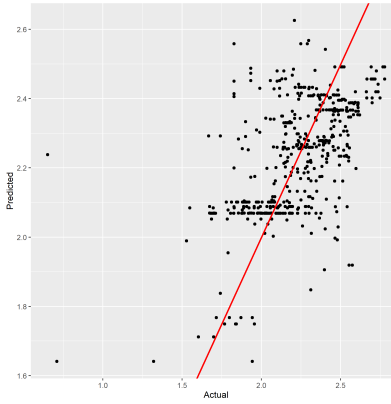


Figure 5: Actual VS Predicted model of *Emilia huxleyi* data. LASSO model.

From the actual vs predicted plots, majority of the points are clustered around the one-to-one line. With an exception of a few outliers, the overall relationships in all three plots show a good fit.

## 6 Discussion and conclusions

In this project, methods using hierarchical generalized additive models as well as LASSO regression were used to build species models in the North Atlantic ocean using R/Rstudio. Initially, the data were compiled using the Coccobase data repository, and subsetted using the longitude and latitude, and further partitioned into the 6 different species specific data groups using the highest number of observations.

The best models selected for the HGAMs were model G using the *Coccolithus* species data, and model GS

with group smoother on salinity, model GI with group smoother on  $N^*$  using the data for *Emiliana huxleyi*. Using these results, the species distribution models can be further interpreted and summarized using the weighted mean and variances of each of the plotted smooths shown in figures 12, 13, and 14 (Appendix). If an individual were to require a more general interpretation, a weighted average of the individual models' predictions can be used, which would imply model averaging rather than model selection.

The best models using LASSO regression was the model using the data for *Emiliana huxleyi*. Although the LASSO coefficients for the *Syracosphaera pulchra* data shrunk two terms, phosphate and latitude to 0, this model was exempted from the model selection by taking into account the sample size for these species. LASSO regression was mainly useful for variable selection. For example, the salinity term for coccolithophore and *Coccolithus* species shrunk to 0, indicating that the term could be removed for future analysis. Although I was expecting the correlated variables such as phosphate, nitrate, and silicate to be shrunk to 0, this was not the case for this analysis, as the models were not particularly sparse. Unfortunately, model selection for LASSO was not based on AIC, but it would be useful to calculate some numerical value for comparisons. In the future, it would be beneficial to determine the variables for prediction using LASSO regression first, then running the HGAMs for faster computational time, as well as determining the differences in species niche.

To improve the models, it may be beneficial to take into account the times of each observations, and build a time-series model of the distribution of species. This would be another interesting approach to determine if or how the species distributions change over time. A sensitivity analysis should be conducted given that some species data (e.g. *Syracosphaera pulchra*) used had a small number of observations in comparison to other species (e.g. *Emiliana huxleyi*). This will be beneficial, as number of observations may influence the overall fit of the models.

## References

- [1] H. Akaike. A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19:716–723, 1974.
- [2] Andrew Gelman. Multilevel (hierarchical) modeling: What it can and cannot do. *Technometrics*, 2006.
- [3] Antoine Guisan, Thomas C Edwards, and Trevor Hastie. Generalized linear and generalized additive models in studies of species distributions: Setting the scene. *Ecological Modelling*, 2002.
- [4] Andrew Irwin and Seoyeon Park. Marine microbial macroecology, 2021.
- [5] Y. Lee and J. A. Nelder. Hierarchical generalized linear models. *Journal of the Royal Statistical Society: Series B (Methodological)*, 1996.
- [6] M. Naud, R. Sheward, A. J. Irwin, and Z. V. Finkel. Coccobase: an updated compilation of coccolithophore field observations. Pre-print, 2020.
- [7] Colleen J. O’Brien, Meike Vogt, and Nicolas Gruber. Global coccolithophore diversity: Drivers and future change. *Progress in Oceanography*, page 27–42, 2016.
- [8] Seoyeon Park. Honours - coccobase.
- [9] Eric J. Pedersen, David L. Miller, Gavin L. Simpson, and Noam Ross. Hierarchical generalized additive models in ecology: An introduction with mgcv. *PeerJ*, 2019.
- [10] Dionysios E. Raitsos, Samantha J. Lavender, Yaswant Pradhan, Toby Tyrrell, Philip C. Reid, and Martin Edwards. Coccolithophore bloom size variation in response to the regional environment of the subarctic north atlantic. *Limnology and Oceanography*, 2006.
- [11] Robert Tibshirani. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society: Series B (Methodological)*, 58:267–288, 1996.

## 7 Appendix

Table 11: Full species list and corresponding number of observations

Species names	Number of observations	Species names	Number of observations	Species names	Number of observations
<i>Acanthoica acanthos</i>	32	<i>Coccolithophore</i>	1064	<i>Helicosphaera wallichii</i>	4
<i>Acanthoica quattropsina</i>	47	<i>Coccolithus pelagicus</i>	572	<i>Helladosphaera cornifera</i>	105
<i>Acanthoica</i> spp.	35	<i>Coccolithus</i> spp.	230	<i>Homozygosphaera spinosa</i>	7
<i>Algirosphaera robusta</i>	58	<i>Coccosphaera</i> spp.	39	<i>Lohmannosphaera</i> spp.	2
<i>Alisphaera gaudii</i>	1	<i>Corisphaera gracilis</i>	8	<i>Michaelsarsia adriaticus</i>	4
<i>Alisphaera unicornis</i>	18	<i>Corisphaera</i> spp.	13	<i>Michaelsarsia elegans</i>	13
<i>Anthosphaera</i> spp.	7	<i>Coronosphaera mediterranea</i>	150	<i>Michaelsarsia</i> spp.	12
<i>Braarudosphaera</i> spp.	1	<i>Cyrtosphaera aculeata</i>	6	<i>Oolithotus fragilis</i>	25
<i>Calcidiscus leptoporus</i>	43	<i>Discosphaera tubifera</i>	11	<i>Ophiaster formosus</i>	12
<i>Calcidiscus leptoporus / Coccolithus pelagicus</i>	1	<i>Emiliana huxleyi</i>	1367	<i>Ophiaster hydroideus</i>	171
<i>Calcidiscus</i> spp.	57	<i>Emiliana huxleyi / Gephyrocapsa</i> spp.	1	<i>Ophiaster</i> spp.	5
<i>Calciopappus caudatus</i>	116	<i>Florisphaera profunda</i>	3	<i>Palusphaera vandellii</i>	9
<i>Calciopappus rigidus</i>	2	<i>Gephyrocapsa ericsonii</i>	19	<i>Pappomonas</i> spp.	1
<i>Calciopappus</i> spp.	10	<i>Gephyrocapsa muelleriae</i>	34	<i>Papposphaera</i> spp.	4
<i>Calciosolenia brasiliensis</i>	3	<i>Gephyrocapsa oceanica</i>	3	<i>Pontosphaera syracusana</i>	3
<i>Calyptrolithina divergens</i>	2	<i>Gephyrocapsa ornata</i>	1	<i>Reticulofenestra parvula</i>	4
<i>Calyptrosphaera sphaeroidea</i>	90	<i>Gephyrocapsa</i> spp.	49	<i>Reticulofenestra</i> spp.	1
<i>Calyptrosphaera</i> spp.	14	<i>Gladiolithus flabellatus</i>	1	<i>Rhabdosphaera clavigera</i>	51
<i>Canistrolithus</i> spp.	1	<i>Helicosphaera carteri</i>	114	<i>Rhabdosphaera hispida</i>	6
<i>Coccolithaceae</i>	281	<i>Helicosphaera</i> spp.	1	<i>Rhabdosphaera</i> spp.	7

Species names	Number of observations	Species names	Number of observations
<i>Rhabdosphaera xiphos</i>	8	<i>Syracosphaera</i> spp.	226
<i>Scyphosphaera apsteinii</i>	5	<i>Syracosphaera tumularis</i>	16
<i>Syracosphaera anthos</i>	12	<i>Thoracosphaera heimii</i>	7
<i>Syracosphaera bannockii</i>	13	<i>Umbellosphaera</i> spp.	7
<i>Syracosphaera borealis</i>	9	<i>Umbellosphaera tenuis</i>	22
<i>Syracosphaera corolla</i>	15	<i>Umbilicosphaera sibogae</i>	32
<i>Syracosphaera coronata</i>	4		
<i>Syracosphaera dentata</i>	18		
<i>Syracosphaera dilatata</i>	14		
<i>Syracosphaera halldalii</i>	15		
<i>Syracosphaera histrica</i>	9		
<i>Syracosphaera marginiporata</i>	18		
<i>Syracosphaera molischii</i>	97		
<i>Syracosphaera nana</i>	1		
<i>Syracosphaera nodosa</i>	11		
<i>Syracosphaera noroitica</i>	5		
<i>Syracosphaera ossa</i>	14		
<i>Syracosphaera prolongata</i>	8		
<i>Syracosphaera pulchra</i>	240		
<i>Syracosphaera rotula</i>	45		

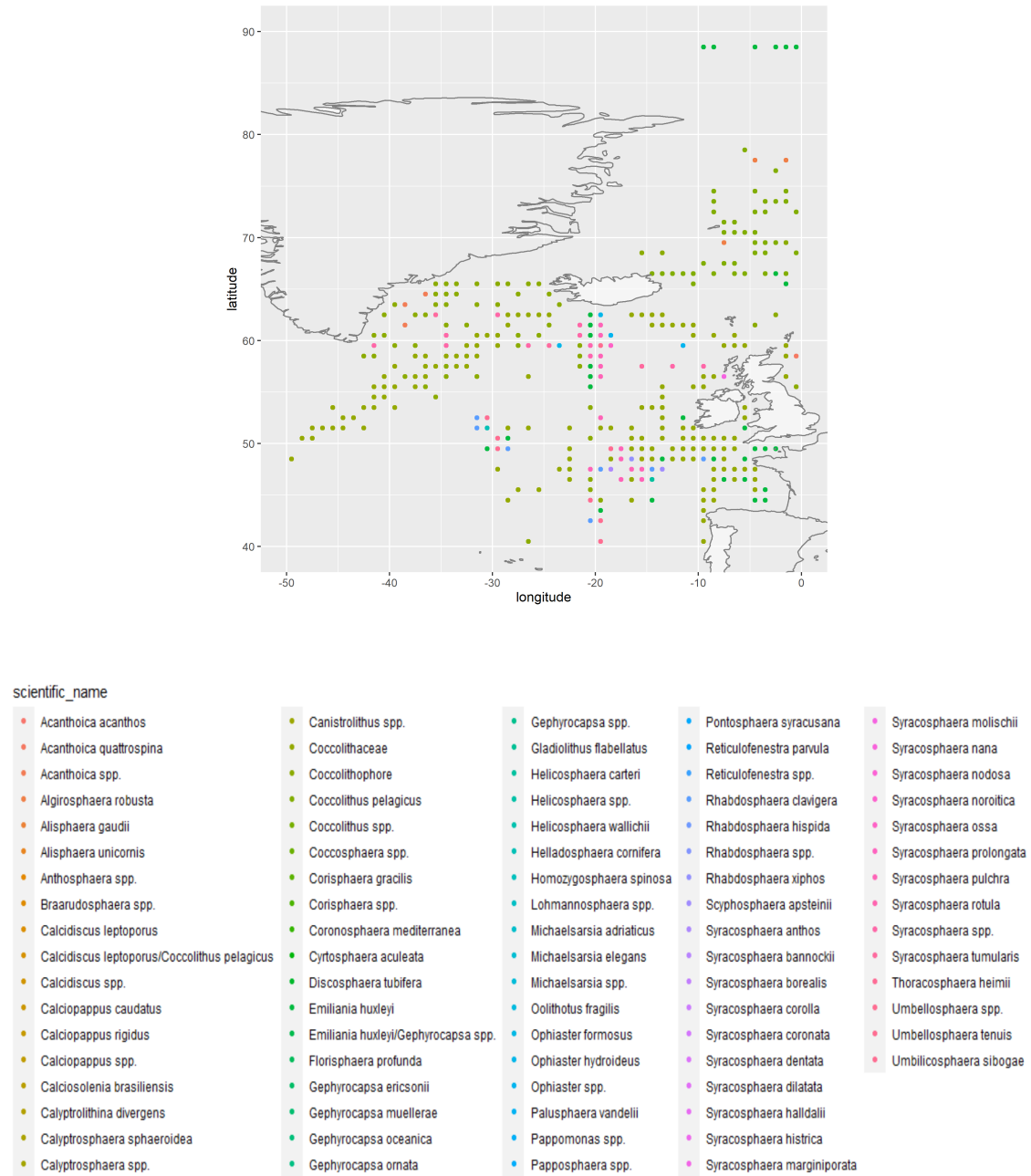


Figure 6: Distribution of all species in full data

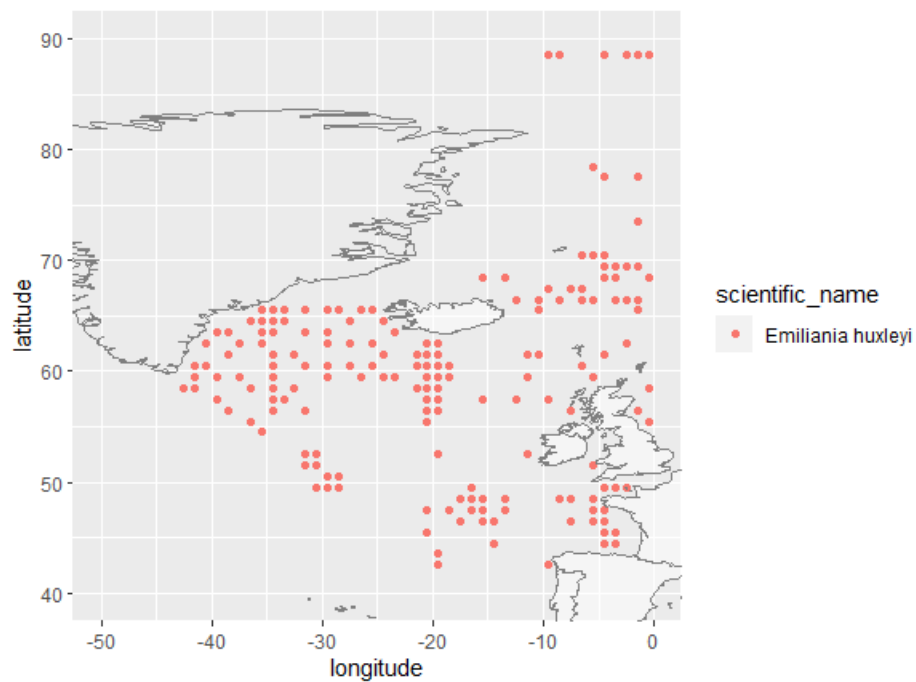


Figure 7: Distribution of *Emiliana huxleyi*

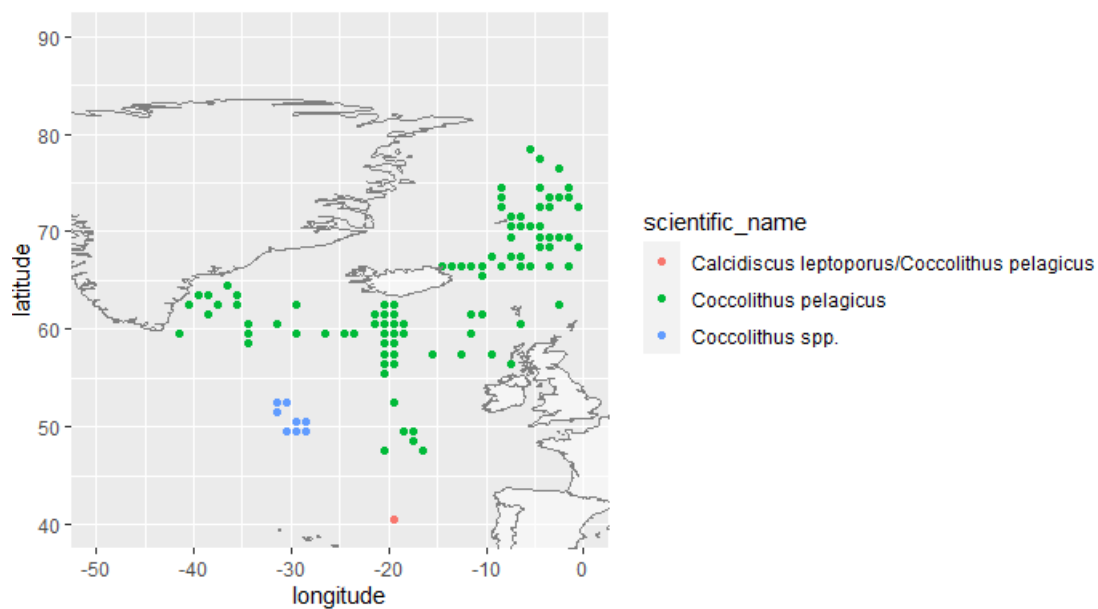


Figure 8: Distribution of *Coccolithus* species

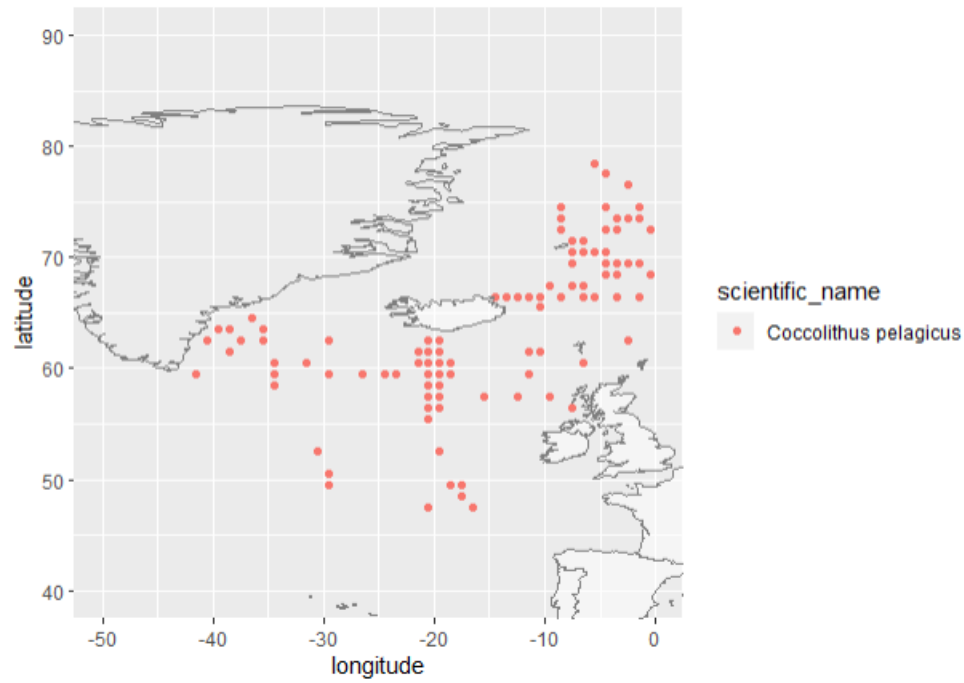


Figure 9: Distribution of *Coccolithus pelagicus*

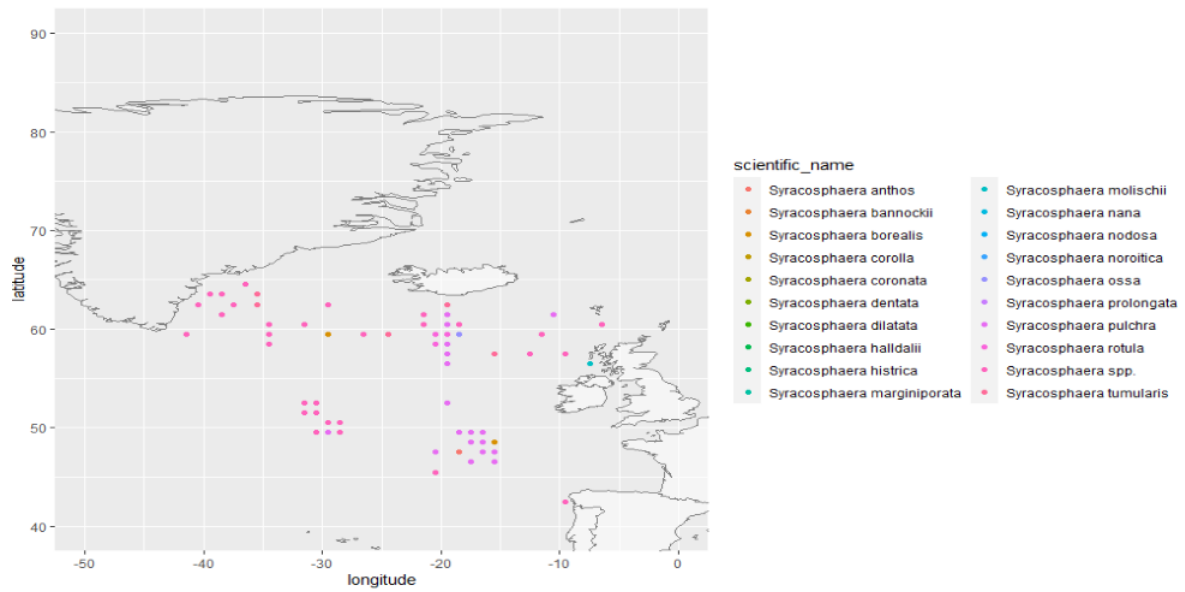


Figure 10: Distribution of *Syracosphaera* species



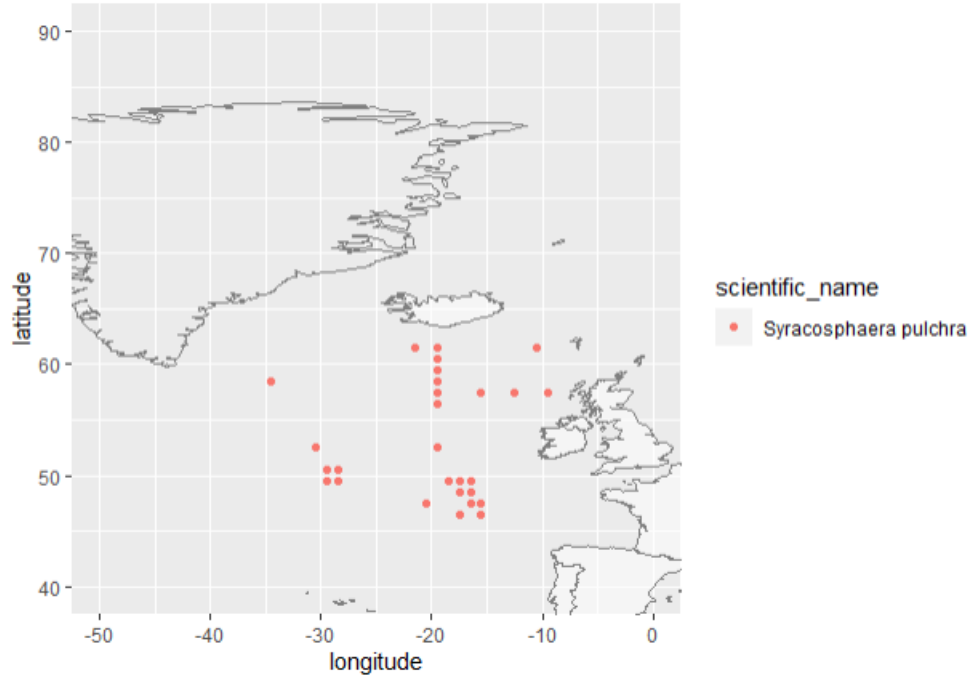


Figure 11: Distribution of *Syracosphaera pulchra*

Next pages will include some GAM plot outputs of training data using "seWithMean = T". These GAM plots correspond to the models that were selected to be the "best" models based on AIC and RMSE values.

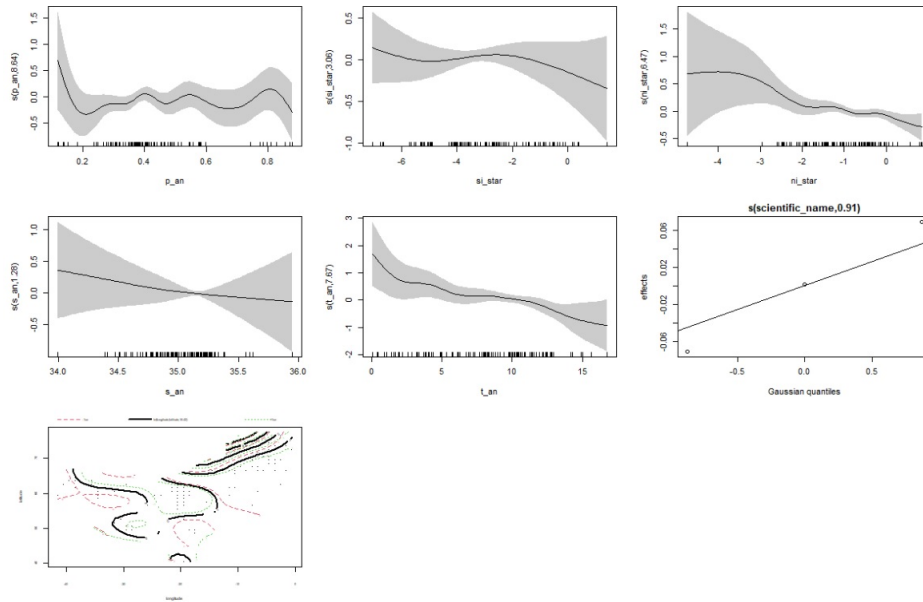


Figure 12: *Coccothithus* species data. Model G

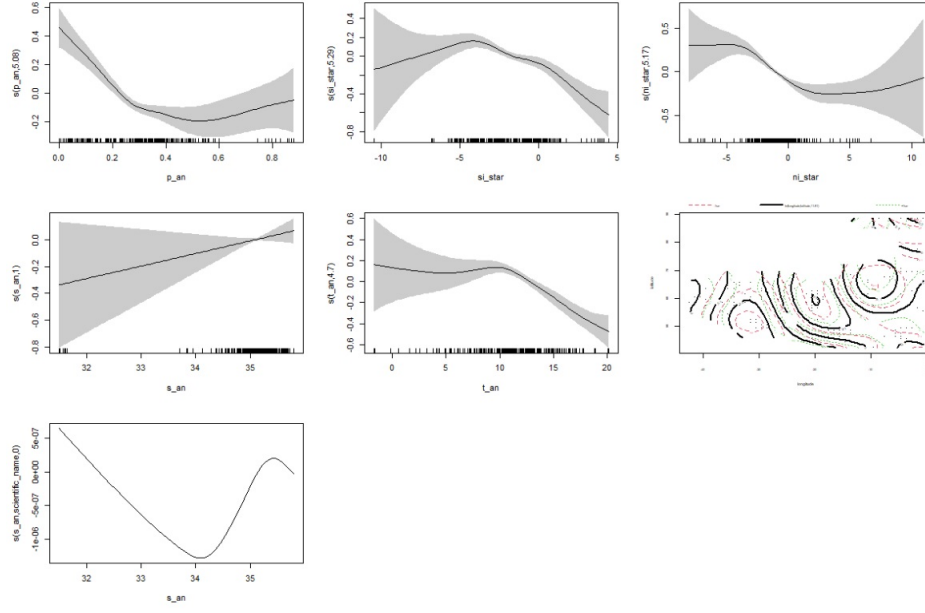


Figure 13: *Emiliana huxleyi* data. Model GS with specific smooth on salinity.

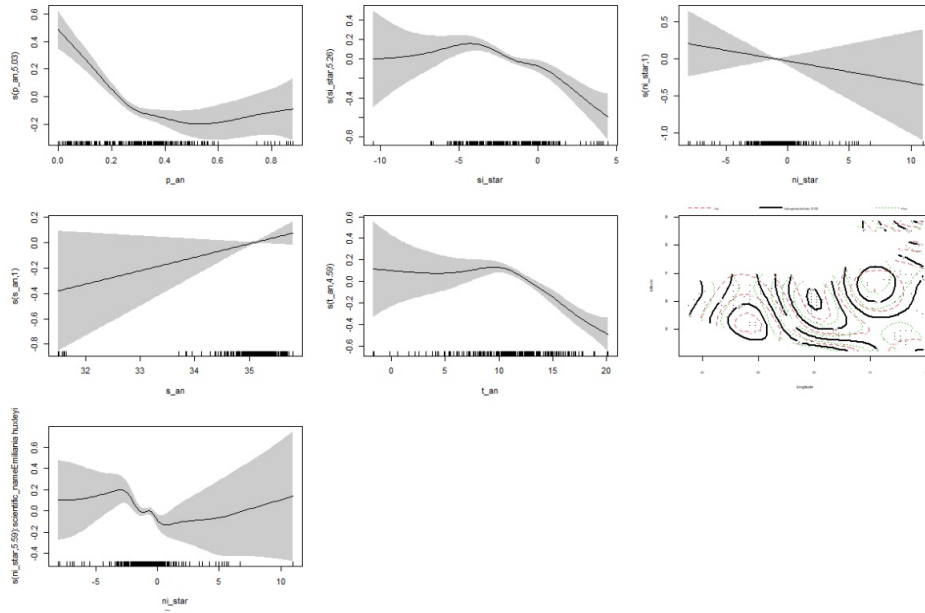


Figure 14: *Emiliana huxleyi* data. Model GI with specific smooth on  $N^*$ .

All code used to create plots and complete analyses can be found at: <https://github.com/calipark1213/Honours-Coccobase>