

Saccharide concentration prediction from proxy ocean samples analyzed via infrared spectroscopy and quantitative machine learning.

Nicole M. North[†], Abigail A. A. Enders[†], Jessica B. Clark[†], Kezia A. Duah[†], Heather C. Allen^{†*}

[†]Department of Chemistry & Biochemistry, The Ohio State University, Columbus, Ohio 43210, United States

Corresponding author

* Heather C. Allen, allen@chemistry.ohio-state.edu

Abstract

Solvated organics in the ocean are present in relatively small concentrations but contribute largely to ocean chemical diversity and complexity. Existing in the ocean as dissolved organic carbon (DOC) and enriched within the sea surface microlayer (SSML), these compounds have large impacts on atmospheric chemistry through their contributions to cloud nucleation, ice formation and other climatological processes. The ability to quantify the concentrations of organics in ocean samples is critical for understanding these marine processes. The work presented herein details an investigation to develop machine learning (ML) methodology utilizing infrared spectroscopy data to accurately estimate saccharide concentrations in complex solutions. We evaluated multivariate linear regression (MLR), K-Nearest-Neighbors (KNN), Decision Trees (DT), Gradient Boosted Regressors (GBR), Multilayer Perceptrons (MLP), and Support Vector Regressors (SVR) toward this goal. SVR models are shown to best predict the accurate generalized saccharide concentrations. Our work presents an application combining fast spectroscopic techniques with ML to analyze organic composition proxy ocean samples. As a result, we target a generalized method for analyzing field marine samples more efficiently, without sacrificing accuracy or precision.

Keywords

sugar, carbohydrate, ocean, spectroscopy, support vector regression, decision trees, gradient boosted regression

Introduction

The sea surface microlayer (SSML) is a multifaceted, deeply complex region of the ocean.¹⁻⁷ As the interface between the Earth's atmosphere and ocean, the SSML performs vital functions that affect climate^{5,8-10} and ice formation.^{4,11-13} Because of unique interfacial anisotropy,¹⁴⁻¹⁷ the physical and chemical properties of the SSML are of interest for their divergence from bulk water behavior. Generally, the SSML is enriched with lipids, proteins, and saccharides (also referred to as sugars or carbohydrates) which contribute to the total dissolved organic carbon (DOC).¹⁸⁻²² Understanding the chemical composition of the SSML provides insight into the biological activity and productivity within the SSML and enables predictions of cloud condensation²³ or ice nucleation,⁴ ultimately aiding climatological models.²⁴⁻²⁷ Recent analyses of saccharide concentrations in SSML have shown concentrations of about 500 nM from eight unique compounds.²⁰ The dynamic nature and chemical complexity of the SSML make monitoring the region difficult, and yet increasingly necessary.

For the described work, glucose and sucrose were chosen as analytes of interest as they are two of the most abundant saccharides found in ocean samples.²⁸ This approach focuses mainly on the quantification of saccharides due to their importance in many marine processes. For example, saccharides are common feedstocks for the ocean ecosystem^{29,30} and can contribute globally to atmospheric processes such as cloud nucleation through transport from the SSML into aerosols.^{28,31} Understanding a generalized saccharide concentration is important to understanding the total ocean chemical diversity and ecosystem health through these processes.

The presented work is motivated by the need for fast, accurate analysis of SSML samples to establish a method that enables exponentially more SSML chemical measurements. Traditional methods to analyze SSML samples are typically limited to mass spectrometry,^{5,32,33} which requires

extensive organic, solid-phase extraction processes. Nevertheless, these methods have provided invaluable information on SSML (and sea spray aerosol) chemical composition. To reduce the sample preparation process and expedite analysis of results, we developed methods that utilize infrared (IR) spectroscopy methods, specifically, attenuated total reflectance Fourier transform infrared (ATR-FTIR) to estimate the saccharide concentration via machine learning (ML) implementations. IR methods provide information on chemical composition and concentration by probing the vibrations of chemical bonds, rather than relying on mass fragmentation. Identification and quantification of specific chemical classes from IR spectra is carried out by analyzing peaks characteristic to specific chemical bonds.³⁴ We note that the limit of detection for ATR-FTIR spectroscopy is higher than for mass spectroscopy, however the speed of analysis for this method is superior.

ML provides a unique avenue to explore relationships among data that cannot be otherwise deduced. The applications to improve or expand chemical systems via ML are broad and present throughout all chemistry fields. Materials design,^{35,36} novel drug discovery,^{37,38} catalyst optimization,^{39,40} and clean energy production^{41,42} are some of the many fields where knowledge has expanded because of ML. Advances in molecular dynamics in combination with machine learning have also paved the way for bridging the connection between molecular structure and physical characteristics.^{43,44} Recent work emphasizes the improved application of FTIR spectroscopy, and more broadly vibrational spectroscopy, for qualitative and quantitative assignment, especially when combined with ML models.^{45,46} Takamura and colleagues explored methods to identify donor biological sex from urine samples.⁴⁷ They presented several ML applications, including partial least-squares discriminant analysis with and without a genetic algorithm, to explore the chemical information contained in their FTIR spectra. They found that

the increased computational complexity of an artificial neural network resulted in comparable results to their discriminant analysis model's predictive power. Butler and coworkers presented successful use of support vector regressors (SVR) in predicting brain cancer from ATR-FTIR spectra.⁴⁸ Their high-throughput approach featured high sensitivity and specificity in the prediction of benign versus malignant samples.

SVRs have also been employed in classification of Raman spectra to identify Alzheimer's Disease in mice; a relevant features map is utilized to identify pertinent peaks that are from molecules known to be associated with the disease. A study from 2022 reports comparable classification accuracy of microplastic Raman microscopy samples from k-nearest neighbors (KNN), multilayer perceptron (MLP), and random forest (RF) models.⁴⁹ These literature examples highlight the diverse applications of ML and develop techniques that expand the applications of chemistry, as we present herein.

This work utilizes ML methods of increasing complexity to evaluate the training data and investigate new data, including field samples with unknown composition. The utilized models in this work are multivariate linear regression, K nearest neighbors, decision trees, gradient boosted regression, multilayer perceptron, and support vector regressors. This diversity in model approach explores the effects of computational complexity, i.e. single models vs ensemble models, and a variety of regression solving techniques.

Fitting data to a linear regression model is common for absorbance data, such as fitting to the Beer-Lambert Law to determine physical constants or identify concentrations of unknown samples.⁵⁰ Absorbance FTIR spectra generally follow a linear relationship of intensity with respect to concentration, which is advantageous for determining new sample composition. Recent work has utilized multiple linear regression to identify heavy metals, including investigating the effect

of surface chemistry on vanadium⁵¹ and lead⁵² toxicity. However, the simplicity of the method ultimately restricts the model's usefulness in more complex, dynamic systems. The largest difference between Beer-Lambert Law linear regression and multivariate linear regression is that all features (in this work, wavenumbers) are used simultaneously to make the multivariate model's assignments.⁵³ This multivariate linear regression (MLR) will act as a benchmark that can be used to compare the other listed models to.

In contrast, SVR fits training data to the best function by minimizing the distance of each value from the fitting equation to be able to predict continuous values. Not all data is appropriate for SVR, but in cases where concentration is being predicted and is linearly correlated with absorbance, it can be a well-suited model. A 2020 report by Mohammadi and colleagues presented an application of SVR to predict different functional group fractions in crude oil.⁵⁴ As another example, ATR-FTIR and SVR were employed by Chen et al. 2022 to predict bio-oil characteristics quickly.⁵⁵

The work described herein provides a discussion on an improved approach to monitoring the SSML. We explore ML approaches to achieve precise and accurate quantitative analysis of simplified proxies of glucose and egg serum albumin (ESA). Glucose is used as our saccharide proxy for training data as it is commonly observed in field measurements and saccharides are frequently reported as a concentration of glucose.^{32,56,57} We also use ESA in our training set because ESA, our SSML protein proxy, has been shown to have surface activity and form insoluble monolayers on aqueous interfaces, despite being a water soluble protein.⁵⁸⁻⁶⁰ While an unlikely protein to find in field samples, ESA provides a complex matrix of amino acids that are abundant in the ocean's water column.^{5,7,61-63} The use of ML in conjunction with vibrational spectroscopy enables greater exploration of chemical space and identifying connections between data. Our

results present, to our knowledge, a first account of predicting saccharide concentration from FTIR spectra of ocean proxy samples using ML.

Methods

Training Solution Preparation, Data Collection, and Data Preprocessing

All chemicals were used as received and all solutions requiring water were prepared using ultrapure water (18 m Ω) from a MilliQ system. For Simplified Proxy (SP) training spectra, stock solutions of 1M glucose (Sigma Aldrich, $\geq 99.5\%$ (GC)) in ultrapure water and 5 mg/mL egg serum albumin (ESA) (Sigma Aldrich, 62-88%, agarose gel electrophoresis) in ultrapure water were prepared. The solution matrix was produced by dispensing the relevant amount of each stock solution via auto pipette and diluting with the requisite amount of water. Briefly, we selected this system and concentrations to have reasonable complexity.

Both the protein and saccharide have IR absorbances from 1800 to 900 cm⁻¹. The peaks were well resolved, with minimal convolution. Inorganic salts were excluded in our matrix, but we provide spectra of the O-H stretching region in the SI to emphasize the limited effect that they have on the IR spectra. Concentrations were selected based on literature precedent from field study results.^{26,27,33} Solutions were measured in triplicate via ATR-FTIR spectroscopy (PerkinElmer Spectrum 3) with a single beam KRS-5/diamond ATR assembly. Spectra were acquired in the “SingleBeam” mode without the use of a continuous reference and were detected using a liquid nitrogen cooled HgCdTe (MCT) detector over 32 scans (approximately one minute) from 4000 to 450 cm⁻¹ with a resolution of 1 cm⁻¹. Spectra were converted to absorbance with a water-only background spectrum (R_o) using the established relationship of $-\log(R/R_o)$. Baseline correction was done using a linear fit model to correct for inconsistent baseline between measurements. Water-only backgrounds were obtained every 5 sample measurements. Triplicate measurements

were used as individual spectra, rather than an average of the three, to provide more machine learning training and testing data (**Figure 1**).

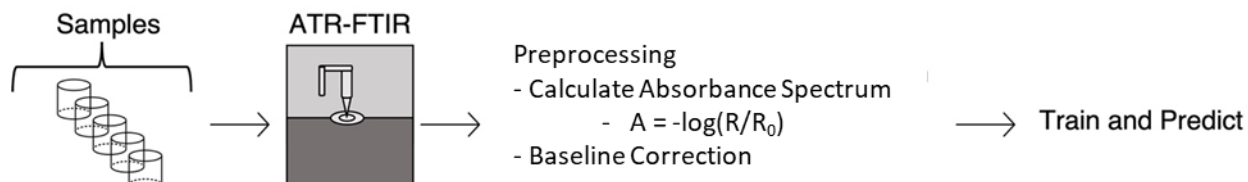


Figure 1. Schematic flow chart of data collection process to the ML pipeline.

Lab Generated Simplified and Ocean Proxy Sample Preparation and Sampling

To test the models' accuracies with increasing chemical complexity, ocean proxy (**OP**) samples were made in the lab with a greater diversity of chemical constituents than the simplified proxies. For these test data, stock ocean proxy-solution was prepared to have 0.1 M sucrose (Sigma Aldrich, $\geq 99.5\%$ (GC)), 0.1 M glucose, 0.5 mg/mL ESA, 3.323 mg/mL bovine serum albumin (BSA) (Sigma Aldrich, $\geq 98\%$, heat shock fraction, pH 7), and 0.1 M 1-butanol (Sigma Aldrich, 99.9%) (**Table 1**). Two additional solutions were prepared via dilution of the stock. The higher concentration dilution was 7.5 mL of stock and 2.5 mL of water and the lower was 5 mL of stock and 5 mL of water. The three solutions were analyzed using the data collection and preprocessing described above.

Table 1. Concentrations of all species in the lab-made ocean proxy samples for evaluation of model accuracy on more chemically diverse conditions

	<i>Ocean Proxy A</i>	<i>Ocean Proxy B</i>	<i>Ocean Proxy C</i>
<i>Concentration of Sucrose (M)</i>	0.10	0.075	0.05
<i>Concentration of Glucose (M)</i>	0.10	0.075	0.05
<i>Concentration of Saccharide (M)</i>	0.20	0.15	0.10
<i>Concentration of ESA (mg/mL)</i>	0.50	0.38	0.25
<i>Concentration of BSA (mg/mL)</i>	3.32	2.49	1.66
<i>Concentration of 1-Butanol (M)</i>	0.10	0.075	0.05

Machine Learning Methods

All machine learning (ML) methods were implemented using Python scripts and SciKit-Learn packages. These are available online at:

https://github.com/Ohio-State-Allen-Lab/Saccharide_Quantification_2024.

Preprocessing

All data, which includes the entire training set of simplified proxy (**SP** – containing only ESA and glucose) and the ocean proxy (**OP** – containing ESA, glucose, BSA, and 1-butanol) samples were standardized using the SciKit-Learn StandardScaler function. This function subtracts the mean of each feature (wavenumber) and divides each feature by the respective standard deviation. The StandardScaler function was first fit using only the SP data, then this fit was applied to both the SP and OP datasets. This was done to avoid the StandardScaler function using the SP dataset information in the OP samples. If the StandardScaler function was fit on the SP and OP datasets together, it would incorrectly inflate the final ability of these models to identify the OP concentrations.⁶⁴ After standardization, the OP data was separated from the data that would then be used for training. The data was then split 70::30 into training and validation/test sets. The latter of which was then split 50::50 into validation and testing datasets. A random state was set to split the data the same way every time into the training, validation, and test datasets to ensure consistency. The training and validation sets were used to train each of the models (210 spectra for training 45 for validation). The withheld test data (45 spectra) were then used to further explore the models' accuracy on previously unseen data that was similar to the data the models were trained on.

A total of 6 machine learning methods were utilized in this work. They will be described here in order of increasing computational complexity.

Multivariate Linear Regression (MLR)

In MLR, all features are fit with a hyperplane in which the dimensionality is determined by the number of features and each feature is has associated weights. This hyperplane is then used to identify concentrations of new samples in the same way that a line would be used for regression with only one feature. Multiple linear models including Lasso, ElasticNet, and Orthogonal Matching Pursuit were tested, but the best performing estimator was the Ridge regressor. This method tends to perform well when there are a large number of features compared to the number of spectral samples.⁶⁵

K-Nearest Neighbors (KNN)

KNN is a method of supervised learning that uses the proximity of previously explored data to make predictions by looking at the distance (the calculation of this distance is variable depending on model parameters) between the neighbors and the training datapoint and using that to adjust the predictions.⁶⁶ In this work, we use the default Minkowski metric for distance which calculates the standard Euclidian distance between points in multivariate space. Different numbers of neighbors between 2 and 10 were tested and the model performed the highest when 5 were used.

Decision Trees (DT)

DTs work to separate the large dataset into smaller pieces repeatedly based on optimized features to be used as split points.⁶⁷ These smallest components, or leaves, then are used to identify predictions for new data. The model utilized in this work terminated splitting once two features were unable to be split further. The model then worked to minimize squared error between training predictions and true values. The original splits were randomized.

Gradient Boosted Regression (GBR)

GBR is an example of an ensemble algorithm that allows for the use of many smaller models, in this context, decision trees.⁶⁸ This method is more computationally complex than a single DT and can identify more complex patterns. The model presented here utilizes a Huber loss function and 2,000 estimators with a learning rate of 0.5 and a max depth of 1.

Multilayer Perceptron (MLP)

MLP is an example of an artificial neural network, a framework of interconnected nodes referred to as neurons.⁶⁹ Each neuron has associated weights, which are adjusted with each training step through a mathematical process of backpropagation. The model presented in this work uses a tanh activation function, an Adam solver, and 500 training steps.

Support Vector Regression (SVR)

SVR utilizes the power of high dimensionality data to identify patterns.⁷⁰ By transforming the data into a higher dimensionality space, it allows for the fitting of the model with different mathematical approaches. The kernel describes the transformation used to transform the data into the high dimensionality hyperplane. This model utilizes a radius bias function (RBF) as the kernel for fitting the dataset.

Model Analysis

To evaluate the models after training, error was also calculated at three different places within the training and testing process. The error calculated is root mean squared error (RMSE). First, the RMSE for the training data is evaluated by comparing the predicted values to the true values with each model. This describes how well the model was able to fit the training data. Next, the validation error was calculated to predict the model's accuracy on new data. Finally, the testing

error focuses on the ability of the model to evaluate data that it has not previously been exposed to.

We also evaluate the prediction of the models on the OP samples to determine how well they perform on data that is chemically different than the data that the models were trained on. This is done by determining the estimation accuracy by comparing the amount of saccharide predicted by each model compared to the true combined saccharide concentration. If the model exactly predicts the concentration, this amount would be 100%. Scores of less than 100% and more than 100% represent under and over prediction respectively. This highlights the degree and directionality of the prediction error in the final estimates of OP data.

Results and Discussion

Evaluating Feasibility of Using IR Spectra to Quantify Saccharide Concentration

The chemical complexity of SP and OP samples is explored with ATR-FTIR spectroscopy and quantitative ML approaches to develop a simple and accurate method of analysis. The FTIR spectra provide chemical information about the sample components and their concentrations, which have a linear correlation with absorbance. The correlation diverges from a linear relationship at high absorbance values, which is not of concern in the presently studied concentration ranges. A single figure containing all the acquired spectra is presented in the SI (Figure S1). Glucose has many vibrational modes that can be used for analysis (Figure S2).

Heat maps can be used to visualize the SP dataset in its entirety. The data was sorted with respect to the concentration of glucose and then plotted against the wavenumber and the intensity at that wavenumber for a given spectrum. This allows for the visualization of the entire dataset in the context of changing glucose concentration and is presented as a heat map in **Figure 2**. A band of increasing intensity can be seen between 1200 and 1000 cm^{-1} correlating to the increasing

concentration of glucose in solution, specifically with the C-C and C-O vibrational modes. The presence of this band supports the ability of the machine learning models to have representative features that will allow for the concentration analysis of glucose.

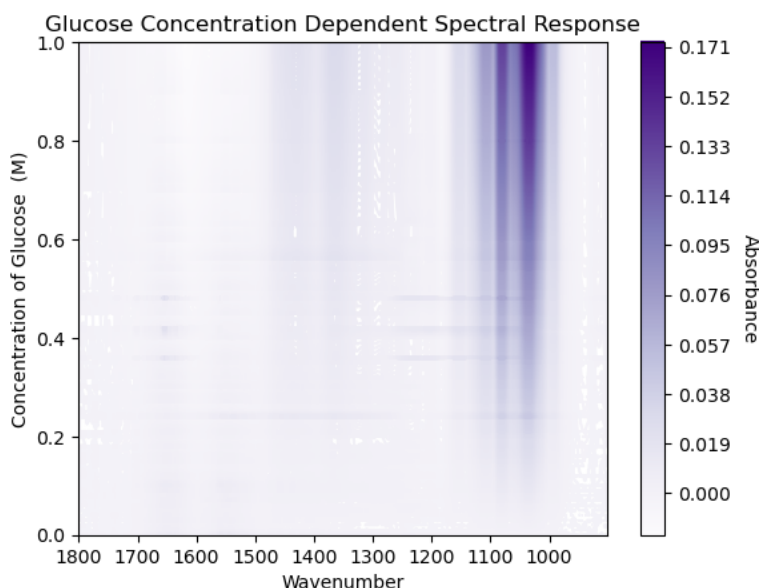


Figure 2. Heat map of the ATR-FTIR dataset as sorted by the concentration of glucose (0 – 1 M). The band of intensity growing in between 1100 and 1000 cm^{-1} corresponds to the increasing C-O stretching within the IR fingerprint region from the increased concentration of glucose. We do not see a strong spectral signature for the ESA relative to that of glucose also in solution (0 – 5 mg/mL) where we would expect the amide bands to exist between 1700 and 1500 cm^{-1} .

To evaluate each model’s ability to accurately predict within the training dataset, model accuracy will be calculated for the training on the simplified proxy (**SP**) dataset. This SP dataset contains only glucose (the analyte of quantification) and ESA (the chemical matrix). To explore if the models are able to expand outside of the explicit training, these models will then be tested on the ocean proxy (**OP**) dataset. Beyond the ESA and glucose within the SP dataset, the OP dataset also contains sucrose, BSA, and 1-butanol. Each of the model’s predicted values will be compared to the additive concentration of glucose and sucrose to make a generalized saccharide concentration (**Figure 3**).

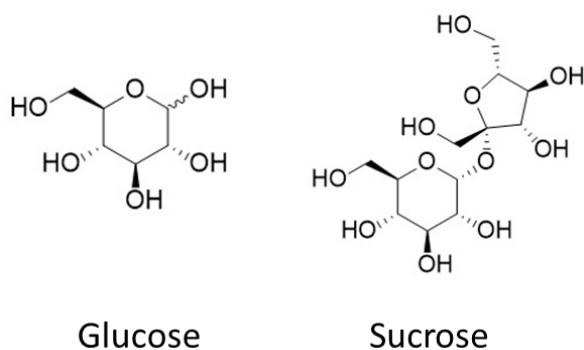


Figure 3. Molecular structures of both glucose (left) and sucrose (right). Both saccharides contain similar vibrational bonds and vibrational environments in regions of the structure. The simplified proxy (**SP**) dataset contains only glucose and egg serum albumin whereas the ocean proxy (**OP**) dataset contains both glucose and sucrose in solution with egg serum albumin, bovine serum albumin, and 1-butanol.

Evaluating Machine Learning Models' Fit of the Simplified Proxy (SP) Dataset

After training, the accuracy of each model's ability to identify the concentrations of the test and validation sets was evaluated to explore the influence of the chosen model to evaluate the SP dataset through analyzing the RMSE error. Ideally, there wouldn't be any effect and the error would be consistent regardless of concentration range. **Figure 4** visualizes these results. DT (**Figure 4 C**) had the smallest associated RMSE and did not exhibit an increased error in low concentrations. KNN, GBR, MLP, and SVR (**Figure 4 B, D, E, and F respectively**) all experienced increased error at low concentrations. R^2 values for each model have also been calculated and are presented in the SI (**Table S2**).

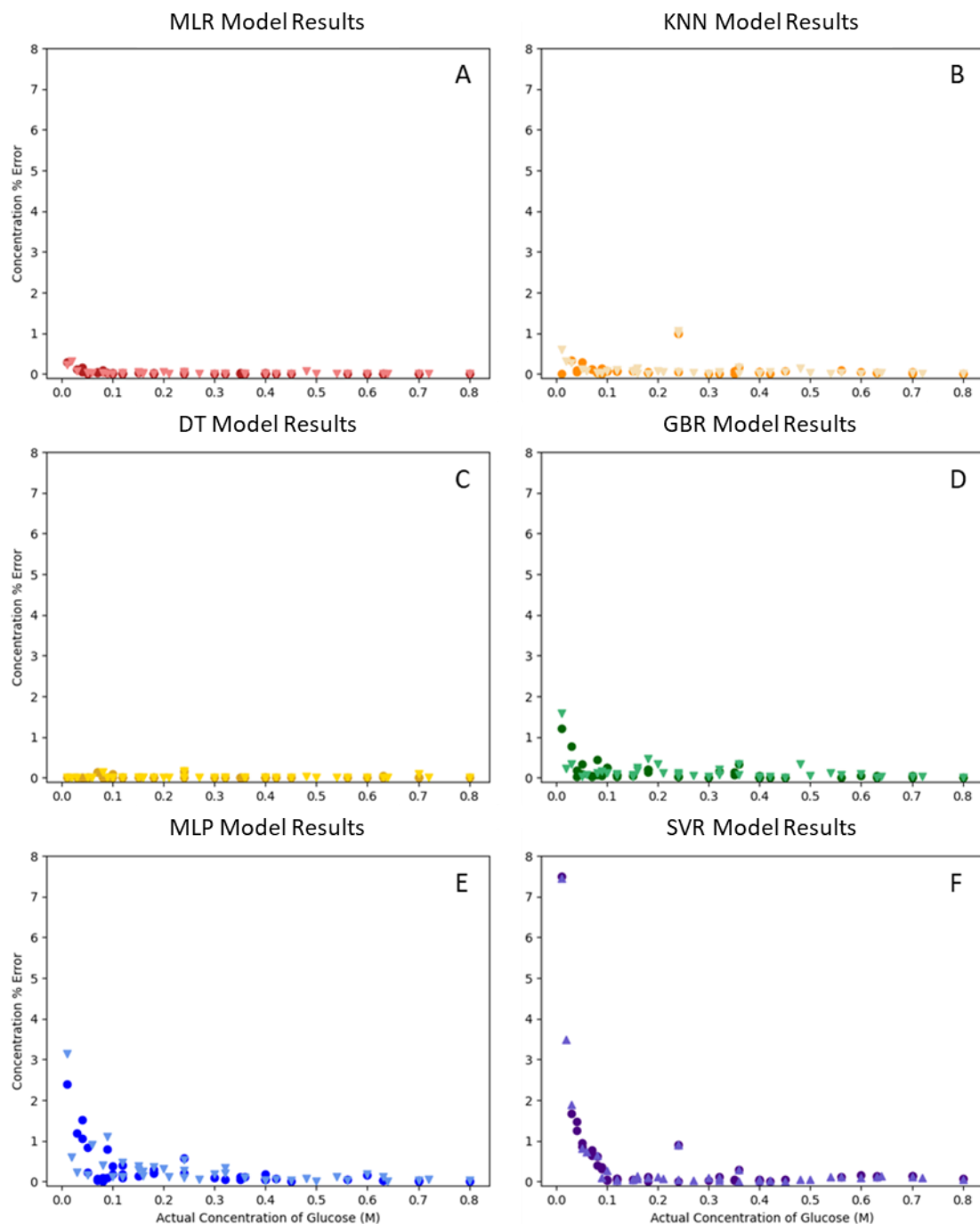


Figure 4(a-f). Scatter plots depicting the accuracy of each of the utilized machine learning models on the simplified proxy (SP) dataset. The y-axis represents the difference between the model assigned and the actual concentrations of the testing dataset divided by the actual concentrations multiplied by 100% (circles) and the withheld validation dataset (triangles). The gradient boosted regression, multilayer perceptron, and support vector regression models do experience an increased error at low concentrations.

To perform a more in-depth error analysis, each of the model's RMSE was calculated between each step of the training by evaluating the training, validation, and test sets' final accuracies. All of the models had smaller than 70 mM in error amongst the different steps. These results have been visualized in **Figure 5**.

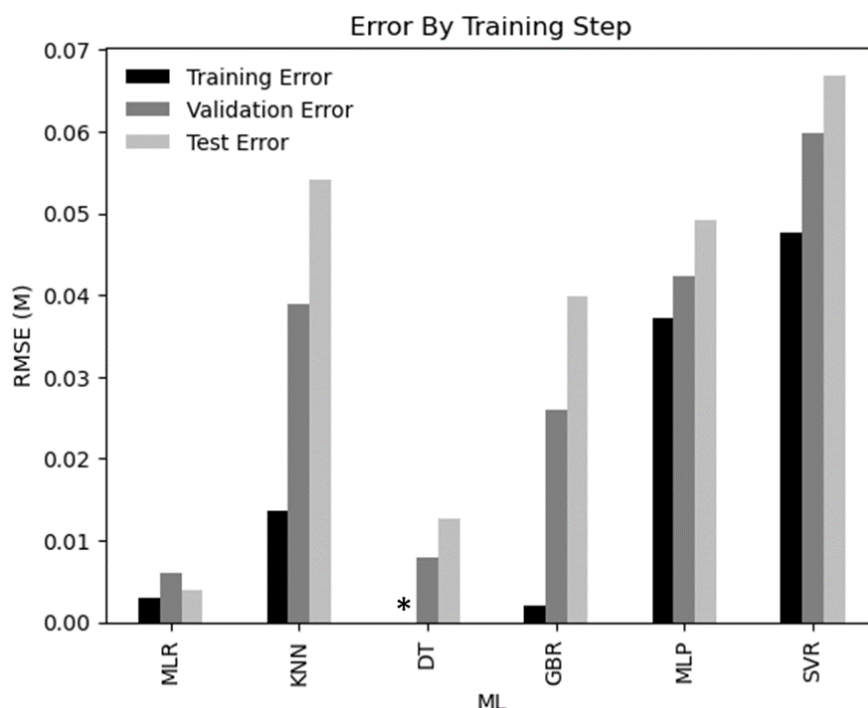


Figure 5. Bar graphs depicting the associated root mean squared error (RMSE) in each part of the training process for the simplified proxy (SP) dataset. All models have a final testing error of less than 0.07 M, but the MLR performed the best in this evaluation. The asterisk indicates that for the decision trees the training error was 0.00 M.

Evaluating Machine Learning Models' Fit of the Ocean Proxy (OP) Dataset

The saccharide concentrations of the OP samples were then estimated using these same ML models. The “true” saccharide concentrations are defined as the sum of the concentrations of glucose and sucrose. This additive concentration, coupled with the increased complexity of the matrix extends these proxies beyond the chemical space that the models were originally trained on. For the purpose of identifying a generalized saccharide concentration, it is important to select

for the models with the highest estimation accuracy when comparing the estimated and true concentrations without disproportionately valuing low or high concentration samples. A model performing poorly here doesn't suggest that the model is poorly trained, just that it doesn't have the capacity to generalize that far beyond the training. For example, MLR had the lowest RMSE error in validation and test datasets as seen in **Figure 5** for the simplified proxies. The MLR, however, only has an estimation accuracy of 50-60% on the OP data, underestimating the combined saccharide concentration by approximately half. This suggests that the MLR model is highly fit to glucose and does not generalize to sucrose, which for other chemical contexts would be ideal.

The highest accuracy in identifying the combined saccharide concentrations came from the SVR and GBR models. They were both able to assign 2/3 of the solutions within 20% of the true concentration of combined saccharide. SVR showed less spread in its predictions but tended to overestimate. The lowest concentration of saccharide was not correctly identified but also existed outside of the range of concentrations where SVR was performing well (**Figure 4**). GBR did not consistently over or underestimate, but it had a large spread in prediction accuracy. These results are shown in **Figure 6**.

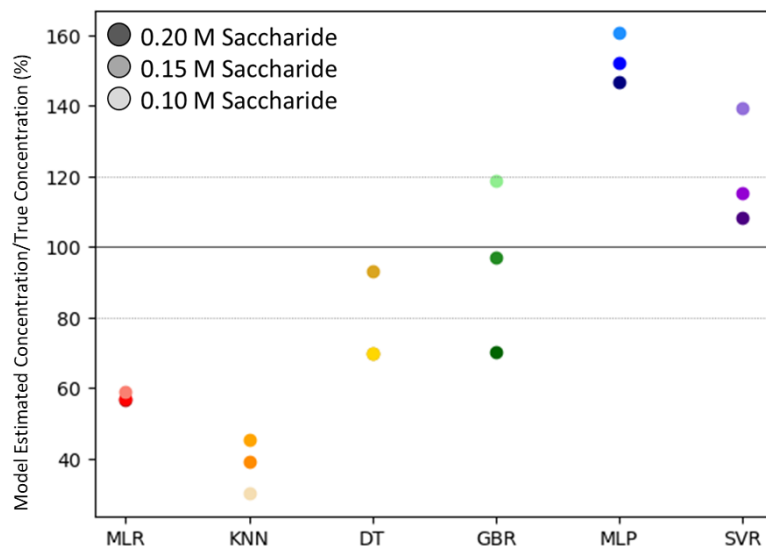


Figure 6. Predicted concentration divided true concentration of combined saccharide for ocean proxy (OP) saccharide concentrations. Solid line at 100 represents 100% meaning that the predicted concentration equals the predicted concentration. The dotted lines represent $\pm 20\%$. The darkest markers in each column represent the highest concentration of saccharide in OP (0.20 M) and the lightest represent the least concentrated (0.10 M). The models have varied levels of success at identifying samples that are far removed from the original training set. The highest performing models were GBR and SVR.

Summary of Discussion

Our quantitative results indicate SVR and GBR are the most promising models to explore for identifying concentrations of saccharides within ocean samples. They are both able to estimate the combined saccharide concentrations within 20% for 2/3 of the complex OP samples. This accuracy would likely be increased if the data that the models were trained on were more chemically similar to the OP dataset as that training would be more relevant to the OP data.

Conclusions

To develop efficient, less-expensive analytical techniques for analysis of the SSML, several ML methods were applied to ATR-FTIR spectra and used to determine saccharide concentration and chemical composition of aqueous samples. Our results indicate that SVR and GBR models are viable for complex solutions, especially considering the training sample data is

relatively simple. The research presented herein provides a unique approach to studying the contributions to the DOC and as a result the SSML utilizing the advanced computational tools available and reduces the time needed to perform analyses of marine samples. Further work should focus on finding an optimal training data set, investigating quantifying other organic concentrations, and intercalating other spectroscopic or spectrometric data, to name a few. An improved understanding and quantification of the marine organics is achievable, wherein more frequent measurements and analysis can occur, ultimately providing more information about the productivity of the marine organics and thus their effects on our atmosphere and climate.

Supplemental Information

Appendix A. ATR-FTIR Spectra of all Training Samples.

Appendix B. Vibrational Analysis of Glucose and ESA

Appendix C. Highest Concentration of ESA and Glucose

Appendix D. Selected Single Peak Beer's Law Analysis

Appendix E. Analysis of MLR and SVR Weights

Appendix F. Tabulated Values for Accuracy and Fit for Each ML Model

Appendix G. Concentration Predictions for Ocean Proxy Solutions for Each ML Model

Acknowledgements

N.M.N. acknowledges funding support from NASA's Future Investigators of NASA Earth and Space Science Technology (FINESST) grant number 20-PLANET20-0067. A.A.A.E. and H.C.A. acknowledge funding from the National Science Foundation through the Center for Aerosol Impacts on Chemistry of the Environment (CAICE) under Grant No. CHE-1801971. J.B.C. acknowledges funding support from the National Science Foundation through Grant No. CHE-2102313.

References

- (1) Carlson, D. J. Dissolved Organic Materials in Surface Microlayers: Temporal and Spatial Variability and Relation to Sea State. *Limnol. Oceanogr.* **1983**, 28 (3), 415–431. <https://doi.org/10.4319/lo.1983.28.3.0415>.
- (2) Cunliffe, M.; Engel, A.; Frka, S.; Gašparović, B.; Guitart, C.; Murrell, J. C.; Salter, M.; Stolle, C.; Upstill-Goddard, R.; Wurl, O. Sea Surface Microlayers: A Unified Physicochemical and Biological Perspective of the Air–Ocean Interface. *Prog. Oceanogr.* **2013**, 109, 104–116. <https://doi.org/10.1016/j.pocean.2012.08.004>.
- (3) Engel, A.; Bange, H. W.; Cunliffe, M.; Burrows, S. M.; Friedrichs, G.; Galgani, L.; Herrmann, H.; Hertkorn, N.; Johnson, M.; Liss, P. S.; Quinn, P. K.; Schartau, M.; Soloviev, A.; Stolle, C.; Upstill-Goddard, R. C.; van Pinxteren, M.; Zäncker, B. The Ocean's Vital Skin: Toward an Integrated Understanding of the Sea Surface Microlayer. *Front. Mar. Sci.* **2017**, 4 (MAY), 1–14. <https://doi.org/10.3389/fmars.2017.00165>.
- (4) Chance, R. J.; Hamilton, J. F.; Carpenter, L. J.; Hackenberg, S. C.; Andrews, S. J.; Wilson, T. W. Water-Soluble Organic Composition of the Arctic Sea Surface Microlayer and Association with Ice Nucleation Ability. *Environ. Sci. Technol.* **2018**, 52 (4), 1817–1826. <https://doi.org/10.1021/acs.est.7b04072>.
- (5) Cochran, R. E.; Laskina, O.; Trueblood, J. V.; Estillore, A. D.; Morris, H. S.; Jayarathne, T.; Sultana, C. M.; Lee, C.; Lin, P.; Laskin, J.; Laskin, A.; Dowling, J. A.; Qin, Z.; Cappa, C. D.; Bertram, T. H.; Tivanski, A. V.; Stone, E. A.; Prather, K. A.; Grassian, V. H. Molecular Diversity of Sea Spray Aerosol Particles: Impact of Ocean Biology on Particle Composition and Hygroscopicity. *Chem* **2017**, 2 (5), 655–667. <https://doi.org/10.1016/j.chempr.2017.03.007>.
- (6) Ault, A. P.; Moffet, R. C.; Baltrusaitis, J.; Collins, D. B.; Ruppel, M. J.; Cuadra-Rodriguez, L. A.; Zhao, D.; Guasco, T. L.; Ebben, C. J.; Geiger, F. M.; Bertram, T. H.; Prather, K. A.; Grassian, V. H. Size-Dependent Changes in Sea Spray Aerosol Composition and Properties with Different Seawater Conditions. *Environ. Sci. Technol.* **2013**, 47 (11), 5603–5612. <https://doi.org/10.1021/es400416g>.
- (7) Bertram, T. H.; Cochran, R. E.; Grassian, V. H.; Stone, E. A. Sea Spray Aerosol Chemical Composition: Elemental and Molecular Mimics for Laboratory Studies of Heterogeneous and Multiphase Reactions. *Chem. Soc. Rev.* **2018**, 47 (7), 2374–2400. <https://doi.org/10.1039/c7cs00008a>.
- (8) Abraham, J. P.; Baringer, M.; Bindoff, N. L.; Boyer, T.; Cheng, L. J.; Church, J. A.; Conroy, J. L.; Domingues, C. M.; Fasullo, J. T.; Gilson, J.; Goni, G.; Good, S. A.; Gorman, J. M.; Gouretski, V.; Ishii, M.; Johnson, G. C.; Kizu, S.; Lyman, J. M.; Macdonald, A. M.; Minkowycz, W. J.; Moffitt, S. E.; Palmer, M. D.; Piola, A. R.; Reseghetti, F.; Schuckmann, K.; Trenberth, K. E.; Velicogna, I.; Willis, J. K. A Review of Global Ocean Temperature Observations: Implications for Ocean Heat Content Estimates and Climate Change. *Rev. Geophys.* **2013**, 51 (3), 450–483. <https://doi.org/10.1002/rog.20022>.
- (9) Burrows, S. M.; Ogunro, O.; Frossard, A. A.; Russell, L. M.; Rasch, P. J.; Elliott, S. M. A Physically Based Framework for Modeling the Organic Fractionation of Sea Spray Aerosol from Bubble Film Langmuir Equilibria. *Atmospheric Chem. Phys.* **2014**, 14 (24), 13601–13629. <https://doi.org/10.5194/acp-14-13601-2014>.
- (10) Cheng, S.; Li, S.; Tsona, N. T.; George, C.; Du, L. Insights into the Headgroup and Chain Length Dependence of Surface Characteristics of Organic-Coated Sea Spray Aerosols. *ACS*

- Earth Space Chem.* **2019**, 3 (4), 571–580.
<https://doi.org/10.1021/acsearthspacechem.8b00212>.
- (11) Wilson, T. W.; Ladino, L. A.; Alpert, P. A.; Breckels, M. N.; Brooks, I. M.; Browse, J.; Burrows, S. M.; Carslaw, K. S.; Huffman, J. A.; Judd, C.; Kilthau, W. P.; Mason, R. H.; McFiggans, G.; Miller, L. A.; Najera, J. J.; Polishchuk, E.; Rae, S.; Schiller, C. L.; Si, M.; Temprado, J. V.; Whale, T. F.; Wong, J. P. S.; Wurl, O.; Yakobi-Hancock, J. D.; Abbatt, J. P. D.; Aller, J. Y.; Bertram, A. K.; Knopf, D. A.; Murray, B. J. A Marine Biogenic Source of Atmospheric Ice-Nucleating Particles. *Nature* **2015**, 525 (7568), 234–238.
<https://doi.org/10.1038/nature14986>.
 - (12) Ting Katty Huang, W.; Ickes, L.; Tegen, I.; Rinaldi, M.; Ceburnis, D.; Lohmann, U. Global Relevance of Marine Organic Aerosol as Ice Nucleating Particles. *Atmospheric Chem. Phys.* **2018**, 18 (15), 11423–11445. <https://doi.org/10.5194/acp-18-11423-2018>.
 - (13) DeMott, P. J.; Hill, T. C. J.; McCluskey, C. S.; Prather, K. A.; Collins, D. B.; Sullivan, R. C.; Ruppel, M. J.; Mason, R. H.; Irish, V. E.; Lee, T.; Hwang, C. Y.; Rhee, T. S.; Snider, J. R.; McMeeking, G. R.; Dhaniyala, S.; Lewis, E. R.; Wentzell, J. J. B.; Abbatt, J.; Lee, C.; Sultana, C. M.; Ault, A. P.; Axson, J. L.; Martinez, M. D.; Venero, I.; Santos-Figueroa, G.; Stokes, M. D.; Deane, G. B.; Mayol-Bracero, O. L.; Grassian, V. H.; Bertram, T. H.; Bertram, A. K.; Moffett, B. F.; Franc, G. D. Sea Spray Aerosol as a Unique Source of Ice Nucleating Particles. *Proc. Natl. Acad. Sci. U. S. A.* **2016**, 113 (21), 5797–5803.
<https://doi.org/10.1073/pnas.1514034112>.
 - (14) Carter-Fenk, K. A.; Dommer, A. C.; Fiamingo, M. E.; Kim, J.; Amaro, R. E.; Allen, H. C. Calcium Bridging Drives Polysaccharide Co-Adsorption to a Proxy Sea Surface Microlayer. *Phys. Chem. Chem. Phys.* **2021**, 23 (30), 16401–16416.
<https://doi.org/10.1039/d1cp01407b>.
 - (15) Yao, X.; Liu, Q.; Wang, B.; Yu, J.; Aristov, M. M.; Shi, C.; Zhang, G. G. Z.; Yu, L. Anisotropic Molecular Organization at a Liquid/Vapor Interface Promotes Crystal Nucleation with Polymorph Selection. *J. Am. Chem. Soc.* **2022**, 144 (26), 11638–11645.
<https://doi.org/10.1021/jacs.2c02623>.
 - (16) Neal, J. F.; Rogers, M. M.; Smeltzer, M. A.; Carter-Fenk, K. A.; Grooms, A. J.; Zerkle, M. M.; Allen, H. C. Sodium Drives Interfacial Equilibria for Semi-Soluble Phosphoric and Phosphonic Acids of Model Sea Spray Aerosol Surfaces. *ACS Earth Space Chem.* **2020**, 4 (9), 1549–1557. <https://doi.org/10.1021/acsearthspacechem.0c00132>.
 - (17) Vazquez De Vasquez, M. G.; Carter-Fenk, K. A.; McCaslin, L. M.; Beasley, E. E.; Clark, J. B.; Allen, H. C. Hydration and Hydrogen Bond Order of Octadecanoic Acid and Octadecanol Films on Water at 21 and 1°C. *J. Phys. Chem. A* **2021**, 125 (46), 10065–10078.
<https://doi.org/10.1021/acs.jpca.1c06101>.
 - (18) Myklestad, S. M. Dissolved Organic Carbon from Phytoplankton. In *Marine Chemistry*; Wangersky, P. J., Ed.; Springer Berlin Heidelberg: Berlin, Heidelberg, 2000; pp 111–148.
https://doi.org/10.1007/10683826_5.
 - (19) Lønborg, C.; Carreira, C.; Jickells, T.; Álvarez-Salgado, X. A. Impacts of Global Change on Ocean Dissolved Organic Carbon (DOC) Cycling. *Front. Mar. Sci.* **2020**, 7 (June), 1–24.
<https://doi.org/10.3389/fmars.2020.00466>.
 - (20) Jayarathne, T.; Sultana, C. M.; Lee, C.; Malfatti, F.; Cox, J. L.; Pendergraft, M. A.; Moore, K. A.; Azam, F.; Tivanski, A. V.; Cappa, C. D.; Bertram, T. H.; Grassian, V. H.; Prather, K. A.; Stone, E. A. Enrichment of Saccharides and Divalent Cations in Sea Spray

- Aerosol during Two Phytoplankton Blooms. *Environ. Sci. Technol.* **2016**, *50* (21), 11511–11520. <https://doi.org/10.1021/acs.est.6b02988>.
- (21) Gericke, A.; Hühnerfuss, H. Investigation of Z- and E-Unsaturated Fatty Acids, Fatty Acid Esters, and Fatty Alcohols at the Air/Water Interface by Infrared Spectroscopy. *Langmuir* **1995**, *11* (1), 225–230. <https://doi.org/10.1021/la00001a039>.
 - (22) Li, Y.; Shrestha, M.; Luo, M.; Sit, I.; Song, M.; Grassian, V. H.; Xiong, W. Salting up of Proteins at the Air/Water Interface. *Langmuir* **2019**, *35* (43), 13815–13820. <https://doi.org/10.1021/acs.langmuir.9b01901>.
 - (23) Orellana, M. V.; Matrai, P. A.; Leck, C.; Rauschenberg, C. D.; Lee, A. M.; Coz, E. Marine Microgels as a Source of Cloud Condensation Nuclei in the High Arctic. *Proc. Natl. Acad. Sci. U. S. A.* **2011**, *108* (33), 13612–13617. <https://doi.org/10.1073/pnas.1102457108>.
 - (24) Ogunro, O. O.; Burrows, S. M.; Elliott, S.; Frossard, A. A.; Hoffman, F.; Letscher, R. T.; Moore, J. K.; Russell, L. M.; Wang, S.; Wingenter, O. W. Global Distribution and Surface Activity of Macromolecules in Offline Simulations of Marine Organic Chemistry. *Biogeochemistry* **2015**, *126* (1–2), 25–56. <https://doi.org/10.1007/s10533-015-0136-x>.
 - (25) Burrows, S. M.; Easter, R.; Liu, X.; Ma, P.-L.; Wang, H.; Elliott, S. M.; Singh, B.; Zhang, K.; Rasch, P. J. OCEANFILMS Sea-Spray Organic Aerosol Emissions – Part 1: Implementation and Impacts on Clouds. *Atmospheric Chem. Phys. Discuss.* **2018**, 1–27. <https://doi.org/10.5194/acp-2018-70>.
 - (26) Elliott, S.; Menzo, Z.; Jayasinghe, A.; Allen, H. C.; Ogunro, O.; Gibson, G.; Hoffman, F.; Wingenter, O. Biogeochemical Equation of State for the Sea-Air Interface. *Atmosphere* **2019**, *10* (5), 1–17. <https://doi.org/10.3390/atmos10050230>.
 - (27) Elliott, S.; Burrows, S.; Cameron-Smith, P.; Hoffman, F.; Hunke, E.; Jeffery, N.; Liu, Y.; Maltrud, M.; Menzo, Z.; Ogunro, O.; Van Roekel, L.; Wang, S.; Brunke, M.; Jin, M.; Letscher, R.; Meskhidze, N.; Russell, L.; Simpson, I.; Stokes, D.; Wingenter, O. Does Marine Surface Tension Have Global Biogeography? Addition for the OCEANFILMS Package. *Atmosphere* **2018**, *9* (6), 216. <https://doi.org/10.3390/atmos9060216>.
 - (28) Hasenecz, E. S.; Kaluarachchi, C. P.; Lee, H. D.; Tivanski, A. V.; Stone, E. A. Saccharide Transfer to Sea Spray Aerosol Enhanced by Surface Activity, Calcium, and Protein Interactions. *ACS Earth Space Chem.* **2019**, *3* (11), 2539–2548. <https://doi.org/10.1021/acsearthspacechem.9b00197>.
 - (29) Liu, N.; Yang, Y.; Li, F.; Ge, F.; Kuang, Y. Importance of Controlling pH-Depended Dissolved Inorganic Carbon to Prevent Algal Bloom Outbreaks. *Bioresour. Technol.* **2016**, *220*, 246–252. <https://doi.org/10.1016/j.biortech.2016.08.059>.
 - (30) Kisand, V.; Tammert, H. Bacterioplankton Strategies for Leucine and Glucose Uptake after a Cyanobacterial Bloom in an Eutrophic Shallow Lake. *Soil Biol. Biochem.* **2000**, *32* (13), 1965–1972. [https://doi.org/10.1016/S0038-0717\(00\)00171-1](https://doi.org/10.1016/S0038-0717(00)00171-1).
 - (31) Jayarathne, T.; Sultana, C. M.; Lee, C.; Malfatti, F.; Cox, J. L.; Pendergraft, M. A.; Moore, K. A.; Azam, F.; Tivanski, A. V.; Cappa, C. D.; Bertram, T. H.; Grassian, V. H.; Prather, K. A.; Stone, E. A. Enrichment of Saccharides and Divalent Cations in Sea Spray Aerosol During Two Phytoplankton Blooms. *Environ. Sci. Technol.* **2016**, *50* (21), 11511–11520. <https://doi.org/10.1021/acs.est.6b02988>.
 - (32) Jayarathne, T.; Sultana, C. M.; Lee, C.; Malfatti, F.; Cox, J. L.; Pendergraft, M. A.; Moore, K. A.; Azam, F.; Tivanski, A. V.; Cappa, C. D.; Bertram, T. H.; Grassian, V. H.; Prather, K. A.; Stone, E. A. Enrichment of Saccharides and Divalent Cations in Sea Spray

- Aerosol During Two Phytoplankton Blooms. *Environ. Sci. Technol.* **2016**, *50* (21), 11511–11520. <https://doi.org/10.1021/acs.est.6b02988>.
- (33) Cochran, R. E.; Laskina, O.; Jayarathne, T.; Laskin, A.; Laskin, J.; Lin, P.; Sultana, C.; Lee, C.; Moore, K. A.; Cappa, C. D.; Bertram, T. H.; Prather, K. A.; Grassian, V. H.; Stone, E. A. Analysis of Organic Anionic Surfactants in Fine and Coarse Fractions of Freshly Emitted Sea Spray Aerosol. *Environ. Sci. Technol.* **2016**, *50* (5), 2477–2486. <https://doi.org/10.1021/acs.est.5b04053>.
- (34) Enders, A. A.; North, N. M.; Fensore, C. M.; Velez-Alvarez, J.; Allen, H. C. Functional Group Identification for FTIR Spectra Using Image-Based Machine Learning Models. *Anal. Chem.* **2021**. <https://doi.org/10.1021/acs.analchem.1c00867>.
- (35) Schleder, G. R.; Acosta, C. M.; Fazzio, A. Exploring Two-Dimensional Materials Thermodynamic Stability via Machine Learning. *ACS Appl. Mater. Interfaces* **2020**, *12* (18), 20149–20157. <https://doi.org/10.1021/acsami.9b14530>.
- (36) Moosavi, S. M.; Jablonka, K. M.; Smit, B. The Role of Machine Learning in the Understanding and Design of Materials. *J. Am. Chem. Soc.* **2020**, *142* (48), 20273–20287. <https://doi.org/10.1021/jacs.0c09105>.
- (37) Batra, K.; Zorn, K. M.; Foil, D. H.; Minerali, E.; Gawriljuk, V. O.; Lane, T. R.; Ekins, S. Quantum Machine Learning Algorithms for Drug Discovery Applications. *J. Chem. Inf. Model.* **2021**, *61* (6), 2641–2647. <https://doi.org/10.1021/acs.jcim.1c00166>.
- (38) Polykovskiy, D.; Zhebrak, A.; Vetrov, D.; Ivanenkov, Y.; Aladinskiy, V.; Mamoshina, P.; Bozdaganyan, M.; Aliper, A.; Zhavoronkov, A.; Kadurin, A. Entangled Conditional Adversarial Autoencoder for de Novo Drug Discovery. *Mol. Pharm.* **2018**, *15* (10), 4398–4405. <https://doi.org/10.1021/acs.molpharmaceut.8b00839>.
- (39) Zhang, J.; Hu, P.; Wang, H. Amorphous Catalysis: Machine Learning Driven High-Throughput Screening of Superior Active Site for Hydrogen Evolution Reaction. *J. Phys. Chem. C* **2020**, *124* (19), 10483–10494. <https://doi.org/10.1021/acs.jpcc.0c00406>.
- (40) Ting, K. W.; Kamakura, H.; Poly, S. S.; Takao, M.; Siddiki, S. M. A. H.; Maeno, Z.; Matsushita, K.; Shimizu, K.; Toyao, T. Catalytic Methylation of M-Xylene, Toluene, and Benzene Using CO₂ and H₂ over TiO₂-Supported Re and Zeolite Catalysts: Machine-Learning-Assisted Catalyst Optimization. *ACS Catal.* **2021**, *11* (9), 5829–5838. <https://doi.org/10.1021/acscatal.0c05661>.
- (41) Miyake, Y.; Saeki, A. Machine Learning-Assisted Development of Organic Solar Cell Materials: Issues, Analyses, and Outlooks. *J. Phys. Chem. Lett.* **2021**, *12* (51), 12391–12401. <https://doi.org/10.1021/acs.jpclett.1c03526>.
- (42) Masood, H.; Toe, C. Y.; Teoh, W. Y.; Sethu, V.; Amal, R. Machine Learning for Accelerated Discovery of Solar Photocatalysts. *ACS Catal.* **2019**, *9* (12), 11774–11787. <https://doi.org/10.1021/acscatal.9b02531>.
- (43) Al Ibrahim, E.; Farooq, A. Transfer Learning Approach to Multitarget Temperature-Dependent Reaction Rate Prediction. *J. Phys. Chem. A* **2022**, *126* (28), 4617–4629. <https://doi.org/10.1021/acs.jpca.2c00713>.
- (44) Freitas, R. S. M.; Lima, A. P. F.; Chen, C.; Rochinha, F. A.; Mira, D.; Jiang, X. Towards Predicting Liquid Fuel Physicochemical Properties Using Molecular Dynamics Guided Machine Learning Models. *Fuel* **2022**, *329*, 125415. <https://doi.org/10.1016/j.fuel.2022.125415>.

- (45) Brandt, J.; Mattsson, K.; Hassellöv, M. Deep Learning for Reconstructing Low-Quality FTIR and Raman Spectra—A Case Study in Microplastic Analyses. *Anal. Chem.* **2021**, *93* (49), 16360–16368. <https://doi.org/10.1021/acs.analchem.1c02618>.
- (46) Fan, X.; Wang, Y.; Yu, C.; Lv, Y.; Zhang, H.; Yang, Q.; Wen, M.; Lu, H.; Zhang, Z. A Universal and Accurate Method for Easily Identifying Components in Raman Spectroscopy Based on Deep Learning. *Anal. Chem.* **2023**. <https://doi.org/10.1021/acs.analchem.2c03853>.
- (47) Takamura, A.; Halamkova, L.; Ozawa, T.; Lednev, I. K. Phenotype Profiling for Forensic Purposes: Determining Donor Sex Based on Fourier Transform Infrared Spectroscopy of Urine Traces. *Anal. Chem.* **2019**, *91* (9), 6288–6295. <https://doi.org/10.1021/acs.analchem.9b01058>.
- (48) Butler, H. J.; Brennan, P. M.; Cameron, J. M.; Finlayson, D.; Hegarty, M. G.; Jenkinson, M. D.; Palmer, D. S.; Smith, B. R.; Baker, M. J. Development of High-Throughput ATR-FTIR Technology for Rapid Triage of Brain Cancer. *Nat. Commun.* **2019**, *10* (1), 1–9. <https://doi.org/10.1038/s41467-019-12527-5>.
- (49) Lei, B.; Bissonnette, J. R.; Hogan, Ú. E.; Bec, A. E.; Feng, X.; Smith, R. D. L. Customizable Machine-Learning Models for Rapid Microplastic Identification Using Raman Microscopy. *Anal. Chem.* **2022**. <https://doi.org/10.1021/acs.analchem.2c02451>.
- (50) Richardson, P. I. C.; Muhamadali, H.; Ellis, D. I.; Goodacre, R. Rapid Quantification of the Adulteration of Fresh Coconut Water by Dilution and Sugars Using Raman Spectroscopy and Chemometrics. *Food Chem.* **2019**, *272* (January 2018), 157–164. <https://doi.org/10.1016/j.foodchem.2018.08.038>.
- (51) Gillio Meina, E.; Niyogi, S.; Liber, K. Multiple Linear Regression Modeling Predicts the Effects of Surface Water Chemistry on Acute Vanadium Toxicity to Model Freshwater Organisms. *Environ. Toxicol. Chem.* **2020**, *39* (9), 1737–1745. <https://doi.org/10.1002/etc.4798>.
- (52) Esbaugh, A. J.; Brix, K. V.; Mager, E. M.; De Schamphelaere, K.; Grosell, M. Multi-Linear Regression Analysis, Preliminary Biotic Ligand Modeling, and Cross Species Comparison of the Effects of Water Chemistry on Chronic Lead Toxicity in Invertebrates. *Comp. Biochem. Physiol. Part C Toxicol. Pharmacol.* **2012**, *155* (2), 423–431. <https://doi.org/10.1016/j.cbpc.2011.11.005>.
- (53) `sklearn.linear_model.LinearRegression`. scikit-learn. https://scikit-learn/stable/modules/generated/sklearn.linear_model.LinearRegression.html (accessed 2023-11-02).
- (54) Mohammadi, M.; Khanmohammadi Khorrami, M.; Vatani, A.; Ghasemzadeh, H.; Vatanparast, H.; Bahramian, A.; Fallah, A. Genetic Algorithm Based Support Vector Machine Regression for Prediction of SARA Analysis in Crude Oil Samples Using ATR-FTIR Spectroscopy. *Spectrochim. Acta. A. Mol. Biomol. Spectrosc.* **2021**, *245*, 118945. <https://doi.org/10.1016/j.saa.2020.118945>.
- (55) Chen, C.; Liang, R.; Ge, Y.; Li, J.; Yan, B.; Cheng, Z.; Tao, J.; Wang, Z.; Li, M.; Chen, G. Fast Characterization of Biomass Pyrolysis Oil via Combination of ATR-FTIR and Machine Learning Models. *Renew. Energy* **2022**, *194*, 220–231. <https://doi.org/10.1016/j.renene.2022.05.097>.
- (56) Schill, S. R.; Burrows, S. M.; Hasenecz, E. S.; Stone, E. A.; Bertram, T. H. The Impact of Divalent Cations on the Enrichment of Soluble Saccharides in Primary Sea Spray Aerosol. *Atmosphere* **2018**, *9* (12), 13–17. <https://doi.org/10.3390/atmos9120476>.

- (57) Roy, S. Distributions of Phytoplankton Carbohydrate, Protein and Lipid in the World Oceans from Satellite Ocean Colour. *ISME J.* **2018**, *12* (6), 1457–1472. <https://doi.org/10.1038/s41396-018-0054-8>.
- (58) Cheng, Y. C.; Bianco, C. L.; Sandler, S. I.; Lenhoff, A. M. Salting-out of Lysozyme and Ovalbumin from Mixtures: Predicting Precipitation Performance from Protein-Protein Interactions. *Ind. Eng. Chem. Res.* **2008**, *47* (15), 5203–5213. <https://doi.org/10.1021/ie071462p>.
- (59) Kudryashova, E. V.; Meinders, M. B. J.; Visser, A. J. W. G.; Van Hoek, A.; De Jongh, H. H. J. Structure and Dynamics of Egg White Ovalbumin Adsorbed at the Air/Water Interface. *Eur. Biophys. J.* **2003**, *32* (6), 553–562. <https://doi.org/10.1007/s00249-003-0301-3>.
- (60) Langmuir, I.; Waugh, D. F. The Adsorption of Proteins at Oil-Water Interfaces and Artificial Protein-Lipoid Membranes. *J. Gen. Physiol.* **1938**, 745–755. <https://doi.org/10.1085/jgp.21.6.745>.
- (61) Angle, K. J.; Nowak, C. M.; Davasam, A.; Dommer, A. C.; Wauer, N. A.; Amaro, R. E.; Grassian, V. H. Amino Acids Are Driven to the Interface by Salts and Acidic Environments. *J. Phys. Chem. Lett.* **2022**, *13* (12), 2824–2829. <https://doi.org/10.1021/acs.jpclett.2c00231>.
- (62) Benner, R.; Kaiser, K. Abundance of Amino Sugars and Peptidoglycan in Marine Particulate and Dissolved Organic Matter. *Limnol. Oceanogr.* **2003**, *48* (1), 118–128. <https://doi.org/10.4319/lo.2003.48.1.0118>.
- (63) Borkowski, M.; Orvalho, S.; Warszyński, P.; Demchuk, O. M.; Jarek, E.; Zawala, J. Experimental and Theoretical Study of Adsorption of Synthesized Amino Acid Core Derived Surfactants at an Air/Water Interface. *Phys. Chem. Chem. Phys.* **2022**, *24* (6), 3854–3864. <https://doi.org/10.1039/D1CP05322A>.
- (64) Lukita, A. *The Dreaded Antagonist: Data Leakage in Machine Learning*. Medium. <https://towardsdatascience.com/the-dreaded-antagonist-data-leakage-in-machine-learning-5f08679852cc> (accessed 2024-02-06).
- (65) *sklearn.linear_model.Ridge*. scikit-learn. https://scikit-learn/stable/modules/generated/sklearn.linear_model.Ridge.html (accessed 2023-12-21).
- (66) *sklearn.neighbors.KNeighborsRegressor*. scikit-learn. <https://scikit-learn/stable/modules/generated/sklearn.neighbors.KNeighborsRegressor.html> (accessed 2023-12-21).
- (67) *sklearn.tree.DecisionTreeRegressor*. scikit-learn. <https://scikit-learn/stable/modules/generated/sklearn.tree.DecisionTreeRegressor.html> (accessed 2023-12-21).
- (68) *sklearn.ensemble.GradientBoostingRegressor*. scikit-learn. <https://scikit-learn/stable/modules/generated/sklearn.ensemble.GradientBoostingRegressor.html> (accessed 2023-12-21).
- (69) *sklearn.neural_network.MLPRegressor*. scikit-learn. https://scikit-learn/stable/modules/generated/sklearn.neural_network.MLPRegressor.html (accessed 2023-12-21).
- (70) *sklearn.svm.SVR*. scikit-learn. <https://scikit-learn/stable/modules/generated/sklearn.svm.SVR.html> (accessed 2023-12-21).

For TOC use only

