

Contents lists available at [ScienceDirect](https://www.sciencedirect.com)

Remote Sensing Applications: Society and Environment

journal homepage: www.elsevier.com/locate/rsase

Exploring the link between spectra, inherent optical properties in the water column, and sea surface temperature and salinity

Solomon White^{a,*}, Encarni Medina Lopez^a, Tiago Silva^b, Evangelos Spyarakos^c,
Adrien Martin^d, Laurent Amoudry^e

^a School of Engineering, University of Edinburgh, Edinburgh, United Kingdom

^b CEFAS, United Kingdom

^c Stirling University, United Kingdom

^d Noveltis, France

^e National Oceanography Centre NOC, United Kingdom

ARTICLE INFO

Keywords:

Salinity
Temperature
Remote sensing
Ocean colour

ABSTRACT

Sea surface salinity and temperature are important measures of ocean health. They provide information about ocean warming, atmospheric interactions, and acidification, with further effects on the global thermohaline circulation and as a consequence the global water cycle. In coastal waters they provide information about sub mesoscale circulations and tidal currents, riverine discharge and upwelling effects. This paper explores the methodology to extract sea surface salinity (SSS) and temperature (SST) from ground based hyperspectral ocean radiance. Water leaving radiance is linked to the inherent optical properties of the water column, effected by the constituent parts. Hyperspectral data at ground level is then used as input to train a linear regression model against temporally and spatially matched water data of SSS and SST. Furthermore, a neural network model to be able to estimate the SST and SSS with the hyperspectral data averaged to multispectral bands to emulate the satellite use case. The neural network model is able to learn the relationship between the multispectral radiance to both SSS and SST values, and can predict these with a root mean square error (RMSE) of 0.2PSU and 0.1 degree respectively. This demonstrates the feasibility of similar algorithms applied to multispectral ocean colour satellites with enhanced coverage and spatial resolution.

1. Introduction

Temperature and salinity are the main determinants of sea water density, influencing ocean circulation from global to local scales, (Jackson et al., 2022) acting on fresh water intrusions around estuaries, MacCreedy and Geyer (2010), and melting icebergs, Giddy et al. (2021). Together they present the most fundamental factors shaping the marine habitat, determining metabolism, and with it survival strategies, Sen Gupta et al. (2020). The shift in global climate, with overall global ocean warming and regional changes is leading to increased thermal and freshwater stratification, (Li et al., 2020), eutrophication, (Breitburg et al., 2018), and hypoxic events, Altieri and Gedan (2015). These changes impact not only the marine ecosystem but also humans through fisheries, tourism and human health, having a particular large impact on coastal economies, Oliver et al. (2018), Minnett et al. (2019). The ability to determine sea surface temperature and salinity enables key measurements in understanding marine heat waves, acidification and de-oxygenation, in response to changing climate, Behrenfeld et al. (2006), Rani et al. (2021), Hughes et al. (2018), Pollock et al. (2014).

* Corresponding author.

E-mail address: solomon.white@ed.ac.uk (S. White).

<https://doi.org/10.1016/j.rsase.2025.101454>

Received 5 September 2024; Received in revised form 19 December 2024; Accepted 7 January 2025

Available online 18 January 2025

2352-9385/Crown Copyright © 2025 Published by Elsevier B.V. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

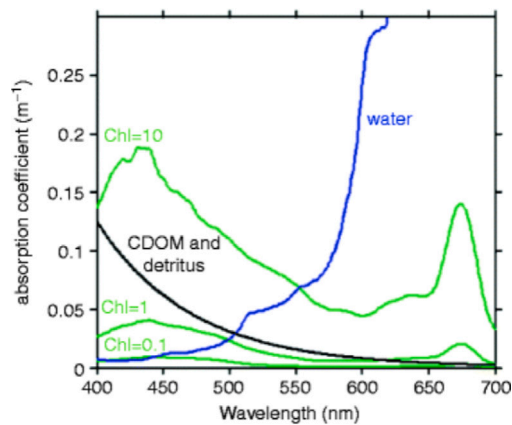


Fig. 1. Absorption coefficient of water constituents as a function of wavelength over the visible spectrum. Bricaud et al. (2004).

The optical properties of the water column are linked to the water leaving radiance, Cael et al. (2020). Inherent Optical Properties (IOPs) such as absorption, backscatter and reflection coefficients at water/air surfaces, can all be estimated from Apparent Optical Properties AOPs (Casey et al., 2020). Fig. 1 shows the absorption coefficients $a(\lambda)$ of water, coloured dissolved organic matter (CDOM) and Chlorophyll, varying with wavelength. By rebuilding the spectral signature of the water leaving radiance $L_w(\lambda)$ in accordance with these water parameters, the relative concentrations can be estimated.

As temperature and salinity only have a negligible effect on the water's optical properties, they must be inferred from the relationship of SST and SSS to the IOPs and AOPs (CDOM, Chl-a, etc.) which do effect the water optical signature and therefore can be measured, Roettgers et al. (2014). Ocean colour algorithms use training data to find the interdependence between, salinity, temperature, Chl-a concentrations and water leaving reflectance ($L_w(\lambda)$). Algorithms are trained on in-situ data from buoys, moorings, drifters or cruises, other remote sensed data (e.g. SMOS, SMAP or satellite estimated data products such as (Moderate Resolution Imaging Spectroradiometer) MODIS SST) and even can be trained using models, O'Reilly and Werdell (2019), Wei et al. (2023).

In the open ocean, salinity change is a conservative process with no source or sink, so the correlation between SSS, Chl, and other freshwater markers can be found empirically and reasonably assumed to be linear, Fournier et al. (2016). With simple regression models finding relationships between the spectral radiance of the freshwater influx, dissolved organic matter (DOM), salinity and temperature (Binding and Bowers, 2003), (Wouthuyzen et al., 2020), in Case 1 conditions, clear oceanic waters that are optically simple and dominated by phytoplankton as the primary source of optical variability (Morel and Prieur, 1977). These simple models suffer in optically complex non conservative Case 2 waters, more varied and turbid coastal or inland waters influenced by various constituents, such as suspended sediments, dissolved organic matter, and phytoplankton, Darecki and Stramski (2004), Hu et al. (2004), Ianson et al. (2003). River discharges, for example, bring large volumes of fresh water. Additionally, sub-mesoscale circulations, such as tidal movements, wind-driven currents, and interactions with underwater topography, all contribute to localized mixing, leading to a non-uniform distribution of salinity as well as vertical stratification. These currents can also lead to sediment resuspension, where particles from the seabed are stirred up by various processes, introducing additional non-conservative elements into the water column, Hsu (2016).

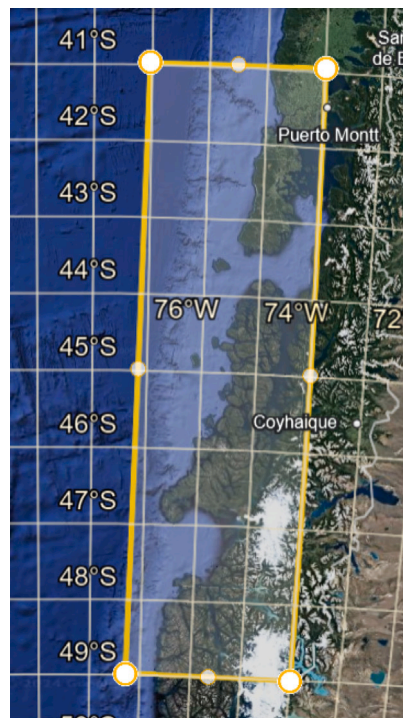
Furthermore, these models are trained and founded on regional coefficients, which are not only area and season specific but vary throughout different water types (e.g. riverine, coral, sea grass, etc.) in the same coastal zone. Therefore, this paper aims to offer a methodology to produce global coverage of SST and SSS, or even regional algorithms with the stability to have accurate estimates, by using a machine learning approach to model design with the offering of higher flexibility.

This paper is structured as follows: (i) introduction to the ground dataset, including hyperspectral and in-water CTD measurements; (ii) development and validation of initial algorithms to determine SST and SSS from optical properties; (iii) exploration of machine learning models to deal with non linearity and produce reliable and scalable observations of SST and SSS.

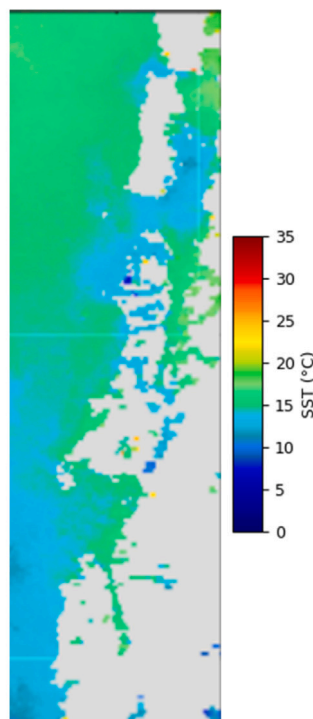
2. Materials and methods

2.1. Patagonia: ground hyperspectral and in-situ SST and SSS data

The study region is located in Patagonia, Chile. The data was collected during a FONDAP-CONICYT (Fund for Research Centres in Priority Areas, by the Chilean Government), funded campaign on-board AGS-61 Cabo de Hornos in Southern Patagonia, from the 15th to the 25th of November 2019. As part of the MEC80180058 and AUB 1900003 Conicyt projects, FONDECYT - National Fund for Scientific and Technological Development (2023). These waters are strongly influenced by terrestrial inputs, providing a unique opportunity to study land-sea interactions and their impact on coastal dynamics. By analysing the spectra and sea water properties in this region, we aim to establish a proof of concept and explore the potential linkages between spectral signatures and



(a) Google Earth view of study region with bounding box.



(b) MODIS Sea surface temperature data

Fig. 2. The Patagonia archipelago region of Chile. Area where the in-situ campaign took place in bounding box. This also shows the two marine terminating glaciers (Strindberg and Ofhidro Glacier) which were sampled during the cruise, at 47° and 49°S respectively. (b) shows SST from MODIS, with the coarse resolution struggles to capture complexities in fjords.

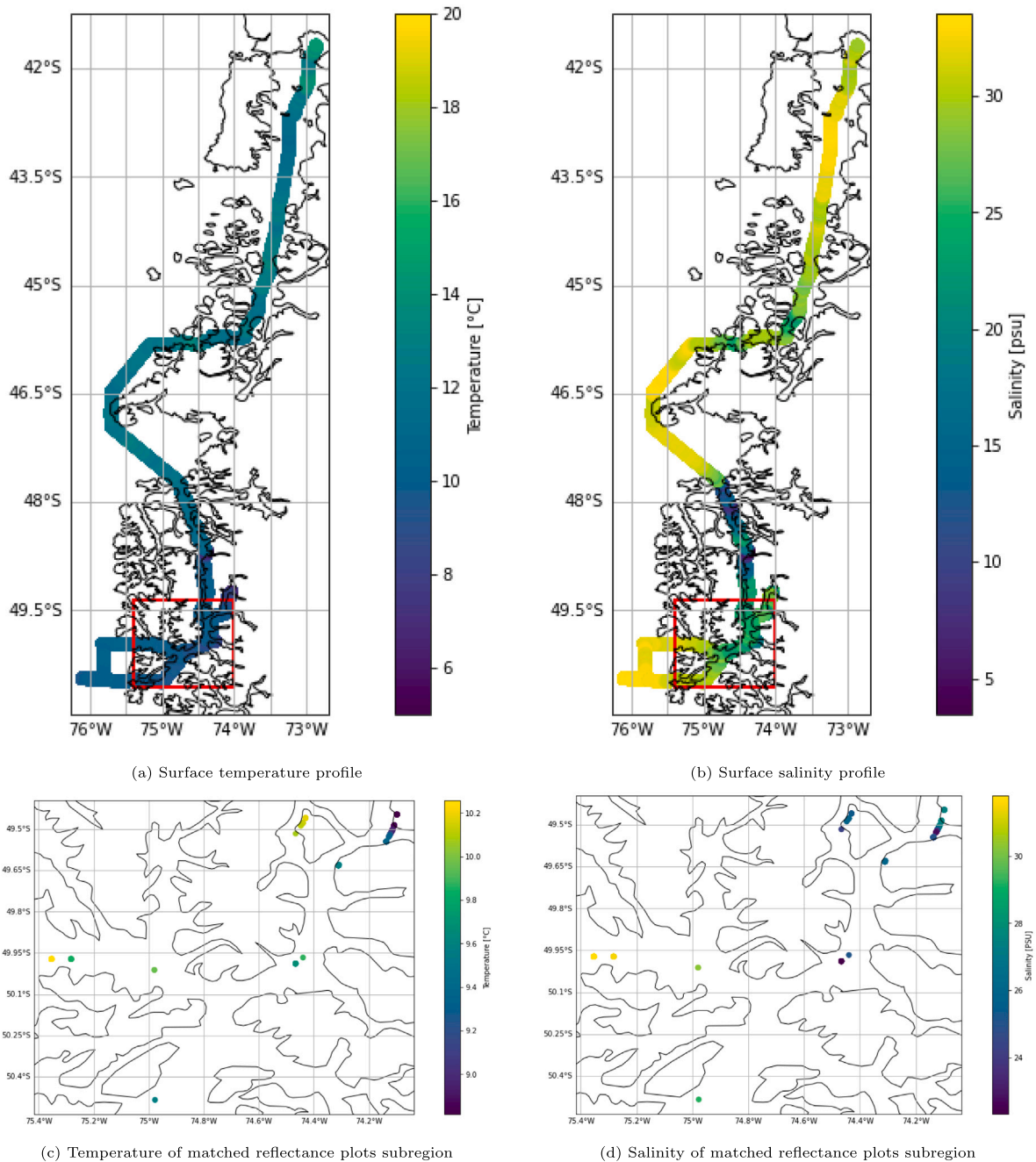
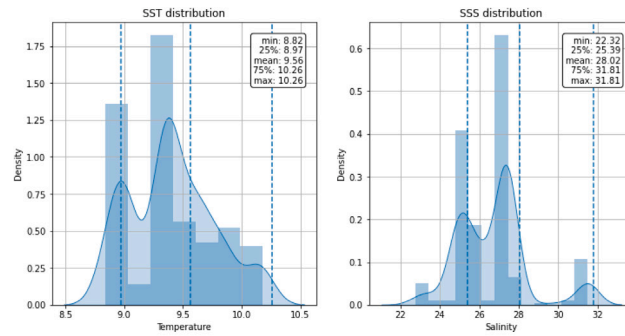
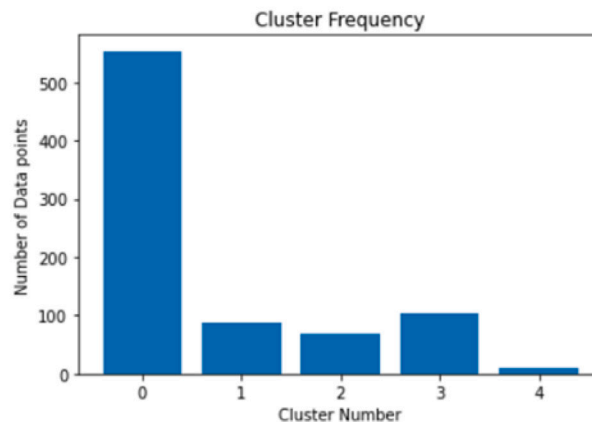


Fig. 3. Temperature (a) and salinity (b) profiles measured by the research vessel and the region of interest (75W, 48S). Influx of cooler fresh glacial water is seen in region 48–49°S, 74 – 75°W, with lower temperature of 5–8 °C compared to average water temperatures of 14–16 °C and salinity values of 10 PSU compared to 30 PSU in surrounding waters. The red box corresponds to the region in (c) and (d) showing the temperature and salinity values for the matched TSG data with reflectance (our model input data). (For interpretation of the references to colour in this figure legend, the reader is referred to the web version of this article.)



(a) Histogram of Temperature (a) and salinity (b) for the matched data sets.



(b) Cluster class frequency plot

Fig. 4. Combined histograms and cluster frequency figures.

environmental parameters. Fig. 2 shows the sampled region, with Fig. 2b showing the MODIS SST imagery for the region, the coarse resolution (1 km² struggles to capture the finer scale temperature changes).

In-situ measurements of water temperature and salinity (through conductivity) every 10 s were taken with a thermosalinograph, mounted at the ship's water intake. Latitude, longitude, time of measurement and depth were also recorded. Hyperspectral in-situ radiometric measurements (324–950 nm, from infrared to ultraviolet electromagnetic spectrum) were collected along with the temperature and salinity data using a ship-mounted set of Trios optical sensors (TriOS (2024)) installed 8 metres above the water surface (with a pixel dispersion of 2.2nm/pixel and wavelength accuracy of 0.2 nm). Remote sensing reflectance R_{rs} was derived from the radiometric measurements using the approach presented by Simis and Olsson (2013). To avoid ship shade and sun glint reflectance measurements were carried out with viewing zenith angle (θ_v) of 40° projecting away from the ship, resulting in a viewing azimuth angle φ^v of 135° away from the solar azimuth φ^s . Relating R_{rs} to the water leaving radiance L_w and fraction p_s of the sky radiance L_s . Spectral filtering using a band pass filter was applied to the reflectance data between the 350–900 nm wavelengths to remove the extreme wavelengths which often suffer from instrument errors. Fig. 3 shows temperature and salinity along the path of the research cruise. Temperature was relatively constant around 10 °C with a cooler region at (75W, 48S) Fig. 3(c), which corresponds to the glacial influx and can be seen in the freshwater signal from the low salinity zone.

The thermosalinograph (TSG) and radiometry data were temporally and spatially matched to each other. Due to the constant monitoring, over 90,000 temperature and salinity measurements were recorded. This number reduced significantly during the matching procedure from the less frequent radiometric (Rrs) values — where the TSG data was temporally matched to the hyperspectral reflectance data. The process explained above reduced the dataset to 490 matched data points.

Fig. 4(a) shows the histogram for temperature and salinity for the matched dataset points. Both display bimodal distributions which seem to correspond to the open and glacial water regions. Temperature has a spread of only 2 °C compared to the 14 °C distribution of all the sampled data, due to the similarity of the reflectance sampling locations. Salinity ranges from more saline open ocean water to the fresher 24PSU due to meltwater influence from the glacier and upwellings from the coastal nature of the cruise.

To further understand the variety and distribution of water types in the dataset, clustering (unsupervised sorting of the data into separate groups) was used with the reflectance values as inputs. A silhouette analysis was used to check the optimal number of

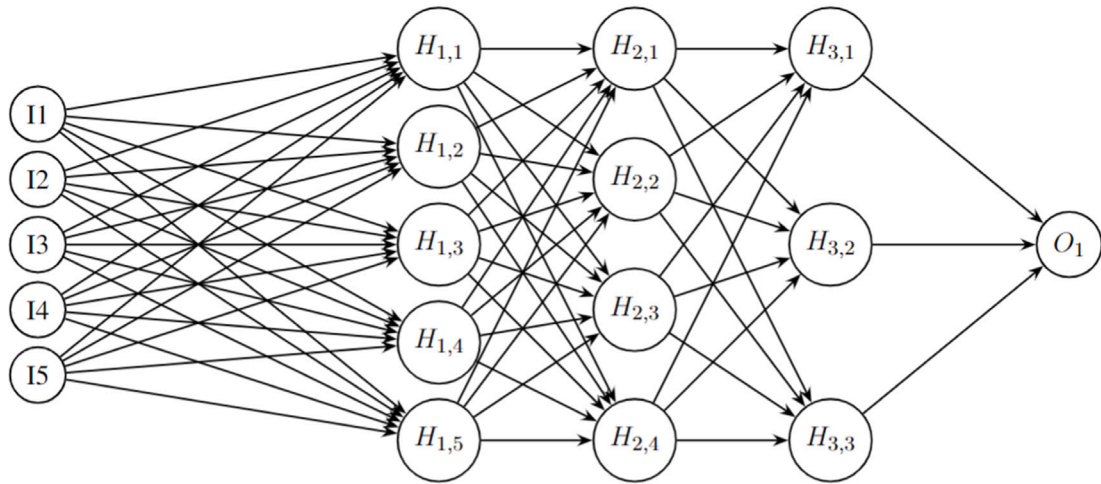


Fig. 5. Neural Network architecture — with 3 hidden layers, 5 inputs and 1 output node. Each individual neuron in each layer is formed by the sum of all the neurons in the previous layer (with separate weights) convolved with some non-linear activation functions.

clusters, scoring each data point based on how well it fit in the assigned class. K-means, spectral and hierarchical clustering, (Society, 1967) all predicted high dominance of a one class, meaning that most of the reflectance values could belong to the same group, with the same water properties (Fig. 4(b)).

2.2. Estimation and regression models

The initial algorithms to estimate the in-situ values are designed to predict the salinity and temperature independently from one another, based purely from the spectral reflectance data. All of the regression and estimation models are the freely available python sci-kit learn libraries, Pedregosa et al. (2011). Firstly, a simple linear regression model is selected with the inputs being the $Rrs(\lambda)$ values (548 bands of 1 nm between 350 and 897 nm) and its output being either temperature or salinity. The model was trained with a train/test split of 0.7/0.3 and root mean squared error (RMSE) and coefficient of determination (R^2 , to capture how much variance is explained by the independent variable) compared for model scoring.

Ridge (Eq. (1)) and Lasso (Eq. (2)). regressions were also tested, both including regularization or penalty terms which act to reduce the number of coefficients, Rajaratnam et al. (2016), Rokem and Kay (2020), (Ranstam and Cook, 2018). The regularization term is controlled by a hyperparameter λ (lambda). As λ increases, the impact of the regularization term grows, leading to a smoother and more generalized model.

$$\min_{\beta} \left\{ \frac{1}{N} \|y - X\beta\|_2^2 + \lambda \|\beta\|_2 \right\} \text{ Ridge regression} \quad (1)$$

$$\min_{\beta} \left\{ \frac{1}{N} \|y - X\beta\|_2^2 + \lambda \|\beta\|_1 \right\} \text{ Lasso Regression} \quad (2)$$

Other models tested were support vector machine regression (SVR), (Ardeshir et al., 2021; López et al., 2022), and decision tree regression, splitting the dataset at each node by minimizing the MSE, James et al. (2013). K-fold cross-validation with 5 splits was also tested to see the model's mean overall error. This approach involves randomly dividing the set of observations into k groups, or folds, of approximately equal size. The first fold is treated as a validation set, and the method is fit on the remaining $k - 1$ folds.

2.3. Neural networks

Neural networks are layered structures linking the input variables via interconnected neurons to an output estimate or class, through a number of hidden layers, as shown in Fig. 5. The advantage of such a structure is that by a combination of these interconnected layers, the network can find non-linear relationships between variables to a high order. Neural networks are excellent for the environmental challenges with large amounts of data and to uncover complex relationships, Hsieh (2009), Domingos (2012). As a final step, a neural network was trained and fitted on the test data to see if the complex model structure was able to better predict the relationship between reflectance and SST/SSS.

A shallow NN with 3 hidden layers and ReLu activation function (a non-linear function between the layers that takes the absolute value of each point) was used, with a sigmoid activation for the final estimate. Finally, to increase training data size and evaluation metrics the dataset was divided into 100 folds using k-fold cross-validation on the 70% of the data designated for training. For each iteration, the model was trained on 99 of the 100 subsets and tested on the remaining fold, repeating this process 100 times

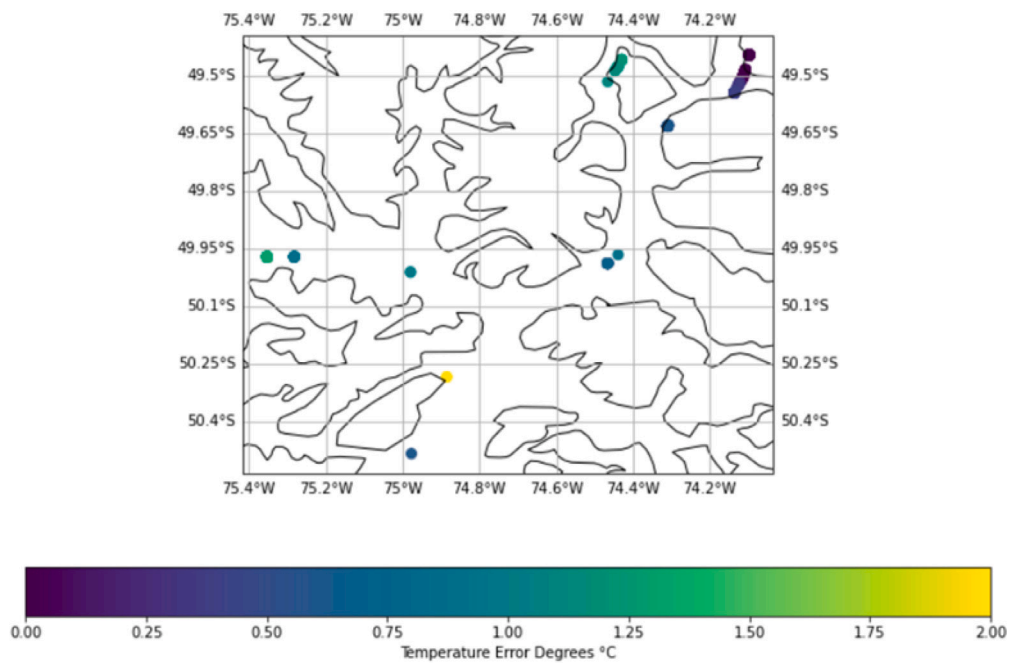


Fig. 6. Distance plot of temperature predicted from PCA (99.9% variance) compared to actual temperature measured.

with a different fold held out for testing each time. This approach ensured that the model is evaluated on different portions of the training data, providing a more robust and generalizable estimate of performance. By aggregating the results from all folds, we obtained average performance metrics, including the standard deviation, which highlights the variability in model performance across different data subsets. The test data was chosen to be a spatially different region to ensure independent evaluation and test the models generalization ability.

3. Results

3.1. Predicting temperature and salinity from hyperspectral radiance — regression results

The linear regression model tends to overfit the data, becoming too related to the training set, which means it is learning relationships from data with no underlying physics. The linear regression model performed very well on the training set, RMSE on the scale of 1^{-10} °C and perfect 1.0 coefficient of determination but when presented with test cases performed comparatively poorly, with polynomial linear regression having RMSE over 1.5 °C. Ridge and Lasso regression perform better with the overfitting less prominent due to the λ regularization term. Table 1 shows the RMSE for these algorithms for both temperature and salinity.

These simple regression models perform well on temperature, with Ridge regression predicting temperature to within 0.5 °C, indeed the polynomial regression clearly shows overfitting. However, salinity performs poorly with simple linear algorithms not able to show the complex relationships present. Polynomial and decision tree give the best RMSE, decision tree being within 1PSU. Salinity requires a more complex model to predict which supports the fact salinity is not a direct indicator to optical properties but instead has to be inferred.

3.2. Feature selection

These models use a large number of inputs with the size of the input array having 548 bands. Typical remote sensing instruments are commonly multispectral, having lower spectral resolution, Vali et al. (2020). Feature selection was used to reduce the dimensionality of the dataset to test the algorithms with lower number of inputs, making it more similar to the multispectral instruments seen on Sentinel-2, Landsat 8 and Sentinel-3.

First, Principal Component Analysis (PCA) was applied to the hyperspectral inputs. PCA is an unsupervised method to reduce the dimensionality of a dataset while maximizing the variance as much as possible, by splitting the data into uncorrelated principle components (eigenvectors). By setting the limit of 99.9% of the variance captured, the data can be reduced to 5 ortho-vectors, with the first and second containing 0.831 and 0.114 of the variance. Fig. 6 shows the spatial error for predicting temperature from the 99.9% variance PCA inputs in the ridge regression model. For capturing 99.5% variance 3 ortho-vectors are needed (see Table 2).

PCA reduces the dimensionality of the dataset by projecting the data in a new multidimensional space, however it is difficult to relate the new vectors to the original bands. Therefore different feature selection methods were tested, keeping the original

Table 1

Temperature and Salinity regression results from hyperspectral data with a 0.7/0.3 train/test split and 490 data points overall. Temperature results are in degrees Celsius (° C) and Salinity results in Practical Salinity Units (PSU). Polynomial Regression is evaluated at power 2 with 15,094 inputs.

Temperature (° C)					
Regression Model	RMSE Train (343)	RMSE Test (147)	MAE Train (343)	MAE Test (147)	K-fold (k=5)
Linear Regression	< 0.01	1.131	< 0.01	1.012	1.312
Ridge	0.469	0.667	0.452	0.601	0.492
Lasso	0.713	0.858	0.690	0.809	0.764
SVM Regression	0.417	0.624	0.400	0.588	0.489
Decision Tree Regression	< 0.01	1.111	< 0.01	1.002	0.613
Polynomial Regression*	< 0.01	1.652	< 0.01	1.520	2.740
Salinity (PSU)					
Regression Model	RMSE Train (343)	RMSE Test (147)	MAE Train (343)	MAE Test (147)	K-fold (k=5)
Linear Regression	< 0.01	3.708	< 0.01	3.641	3.845
Ridge	2.227	2.327	2.180	2.290	2.248
Lasso	3.014	2.975	2.950	2.900	3.003
SVM Regression	1.840	1.915	1.790	1.870	1.808
Decision Tree Regression	1.957	1.900	1.890	1.850	0.899
Polynomial Regression*	< 0.01	1.569	< 0.01	1.540	1.627

Table 2

Temperature and Salinity Regression Tables for the selected input values. Showing the number of inputs for each feature selection, the input bands spectral coverage, the resulting test RMSE, and test MAE.

Temperature (° C) — ridge regression				
Name	Number inputs	Input band coverage	Test RMSE (° C)	Test MAE (° C)
PCA 99.9% variance	5	n/a	0.638	0.532
PCA 99.5%	3	n/a	0.673	0.581
Variance Thresholding (0.0001)	104	489nm–592 nm	0.668	0.563
Top Regression coefficient	10	572–581 nm	0.768 (tree = 0.643)	0.651
10 nm wavelength averaged	10	409–899 nm	0.624	0.526
S2 Multispectral averaged	8	443–865 nm (not possible SWIR)	0.891	0.773
Salinity (PSU) — decision tree regression				
Name	Number inputs	Input band coverage	Test RMSE (PSU)	Test MAE (PSU)
PCA 99.9% variance	5	n/a	0.830	0.703
PCA 90%	2	n/a	1.443	1.193
Variance Thresholding	104	489nm–592 nm	1.127	0.951
Top 10 Regression coefficient	10	554–563 nm	0.920	0.771
10 nm wavelength averaged	10	409–899 nm	1.044	0.876
S2 Multispectral averaged	8	443–865 nm (not possible SWIR)	1.244	1.045

	Rrs489	Rrs490	Rrs491	Rrs492	Rrs493	Rrs494	Rrs495	Rrs496	Rrs497	Rrs498	...	Rrs583	Rrs584	Rrs585	Rrs586	Rrs587
0	0.024075	0.024127	0.024179	0.024248	0.024342	0.024436	0.024524	0.024554	0.024584	0.024613	...	0.025782	0.025664	0.025585	0.025505	0.025426
1	0.014769	0.014860	0.014950	0.015046	0.015151	0.015255	0.015355	0.015418	0.015480	0.015543	...	0.014743	0.014534	0.014304	0.014074	0.013845
2	0.016241	0.016338	0.016434	0.016536	0.016646	0.016756	0.016861	0.016925	0.016988	0.017052	...	0.016092	0.015869	0.015618	0.015366	0.015115
3	0.005057	0.005106	0.005154	0.005202	0.005250	0.005299	0.005345	0.005375	0.005405	0.005435	...	0.005753	0.005650	0.005534	0.005417	0.005301
4	0.001955	0.001957	0.001959	0.001967	0.001984	0.002000	0.002015	0.002018	0.002021	0.002024	...	0.001570	0.001526	0.001487	0.001448	0.001409
...
485	0.010468	0.010454	0.010441	0.010455	0.010509	0.010563	0.010613	0.010628	0.010644	0.010659	...	0.010000	0.009927	0.009849	0.009770	0.009691
486	0.010013	0.010008	0.010004	0.010028	0.010097	0.010166	0.010229	0.010237	0.010245	0.010253	...	0.009838	0.009765	0.009693	0.009621	0.009549
487	0.009892	0.009882	0.009872	0.009891	0.009953	0.010015	0.010072	0.010091	0.010109	0.010128	...	0.009548	0.009477	0.009405	0.009332	0.009259
488	0.008607	0.008606	0.008605	0.008631	0.008695	0.008760	0.008820	0.008839	0.008859	0.008878	...	0.008527	0.008450	0.008381	0.008312	0.008242
489	0.007080	0.007090	0.007101	0.007124	0.007168	0.007211	0.007252	0.007273	0.007294	0.007314	...	0.006930	0.006858	0.006778	0.006699	0.006620

Fig. 7. Feature selection using Variance Threshold 104 bands selected. Columns are the 1 nm hyperspectral bands 489–592 nm and each row is a separate matched data point. Each cell contains the recorded Rrs value for that wavelength at the particular location.

coordinate system. Variance Thresholding was tested, keeping all bands that contribute to the variance of the inputs above 0.0001 reflectance. This resulted in selecting 104 bands between 489–592 nm (corresponding to the blue to yellow part of the electromagnetic spectrum), which and the inputs which provide the most meaningful information for the model predictions. Fig. 7 shows the new inputs following this thresholding.

However, this approach also struggles to capture useful information about the data because the bands are not independent variables (high correlation between neighbouring bands) and each one contains only a small fraction of the variance. This can be seen by the mutual independence coefficient against wavelength showing high similarity between all bands (Fig. 8).

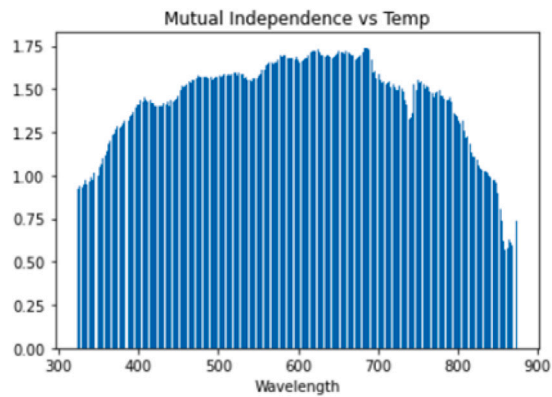


Fig. 8. Mutual independence vs wavelength. The mutual independence coefficient is a dimensionless measure of information about one variable contained in another. It takes values of 0 if the information is independent and the higher the value, the more closely linked variables are, Cover and Thomas (2006).

Univariate feature selection was then used to relate the regression coefficient of each input wavelength, in order to see which spectral bands give the largest weighting to the prediction values.

Both salinity and temperature are most closely related to those bands corresponding to the blue/green colour wavelengths (salinity has been inversely correlated to Chl-a and CDOM which effect the green bands), with temperature also being effected by NIR and red bands. However, selecting the top 10 1 nm hyperspectral bands purely based on this regression coefficient for either salinity or temperature will result in all 10 inputs being in the same narrow wavelength band, all neighbouring. This approach may capture the information for this case study but will have poor diversification to other optical water types, and cannot relate to general multispectral data. To avoid this issue, the bands were averaged into 10 nm bands, then the top 10 were selected within 60 nm spectral segments over the whole frequency range, to ensure a broad input band coverage.

In the final case, the hyperspectral bands were averaged into the multispectral bands matching Sentinel-2, the ESA multispectral satellite with 10 m spatial resolution, in order to relate hyper and multispectral and demonstrate the ability of this process to proceed to satellite top of atmosphere data. Emulated Sentinel 2 data was produced by averaging the hyperspectral 1 nm data into the same bandwidths as the S2 multispectral bands, (European Space Agency, 2019), it results in 8 averaged bands instead of the 13 Sentinel 2 due to lack of the hyperspectral (short wave infrared) SWIR spectral coverage. Fig. 9 shows the regression coefficient weights for the emulated Sentinel 2 data. Showing the coefficient between each general band (Coastal, blue, NIR...) relating to estimation of SST/SSS.

Table 2 shows the performance of each of these condensed input selections. All the inputs for temperature were performed with the ridge regression as that produced the best results for both test RMSE, MAE and K-fold RMSE when implemented on all the hyperspectral inputs, whereas salinity was tested with decision tree regression which performed the best on those metrics.

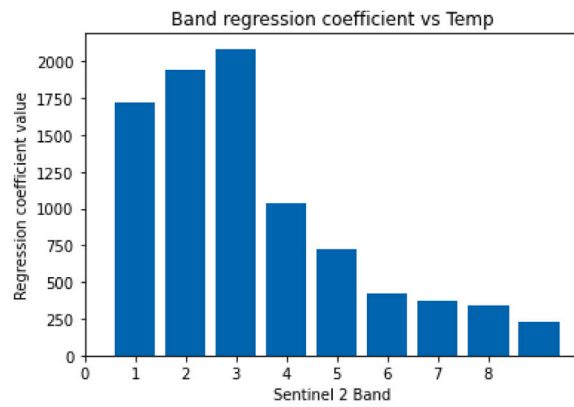
3.3. Neural network implementation

As the number of input spectral bands is reduced, using a more complex model, like a neural network, can better describe the non-linear relationships between colour and surface properties. This performs as well with 10 bands from the 10 nm averaged regression coefficient over all frequencies, as the simple regression models which had 550 input bands.

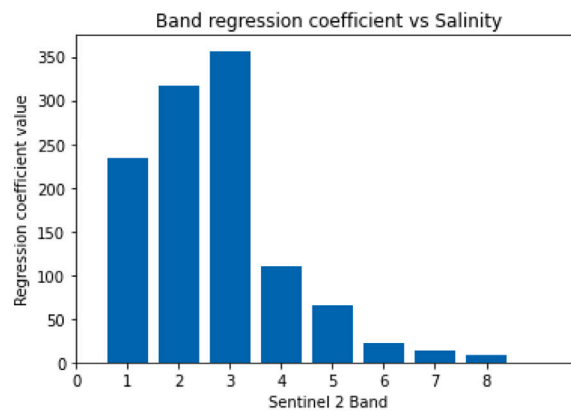
Fig. 10 shows the training and validation loss decreasing with the number of training epochs in the model fitting process. Salinity reduces to a training and validation loss of 0.234 MSE in early epochs and stops after equalizing. The test results show a RMSE of 0.406 PSU, a significant improvement from previous approaches. Temperature is predicted with a validation loss of 0.119, and predicts the temperature for test set with RMSE 0.354 °C. Against the 100 K-fold split the Neural network performance improved. For temperature the mean train RMSE is 0.09 °C with mean validation RMSE of 0.1408 °C \pm 0.0672 std, which shows some variance in model performance on the validation set. The Test RMSE for temperature on the independent data was 0.167 °C. Salinity also showed improved results with training RMSE 0.006PSU, validation RMSE of 0.3271PSU \pm 0.324 (a notably higher std) and test RMSE of 0.391PSU.

4. Discussion

This paper has presented a methodology for extracting salinity and temperature independently from each other using hyperspectral water reflectance values. This shows the capability of ocean colour sensors to infer these physical oceanographic variables. Different regression methods found estimates for salinity and temperature from all the hyperspectral bands. More complex models like decision tree, SVM based regression, and polynomial regression, were able to capture the non-linear relationship with salinity and temperature, with accuracies of 0.899 PSU and 0.492 °C respectively. Through feature selection of pertinent bands, neural network algorithms were able to provide accurate estimates within 0.406 PSU for salinity and 0.354 °C for temperature in test data.

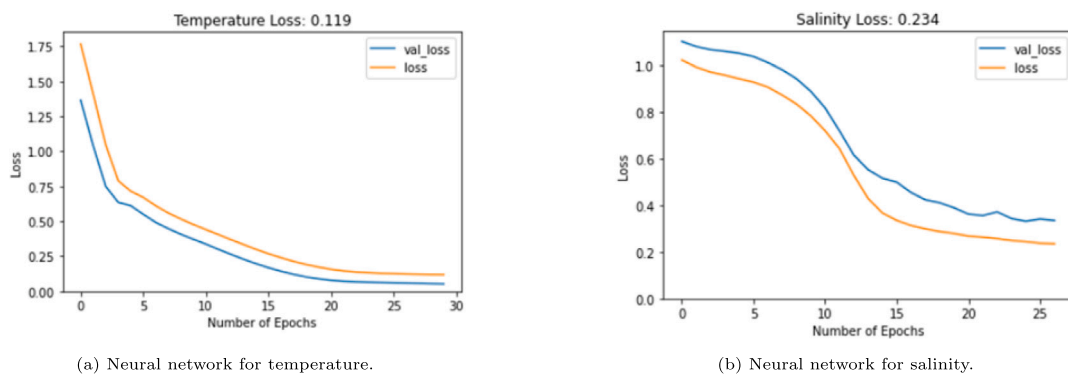


(a) S2 regression coefficients for temperature



(b) S2 regression coefficients for salinity.

Fig. 9. Sentinel 2 regression coefficient values for the predictive model against each band for temperature (a) and salinity (b). A higher coefficient values corresponds to a larger importance for that band in the model prediction.



(a) Neural network for temperature.

(b) Neural network for salinity.

Fig. 10. Neural network train and validation loss (Mean Squared Error) for temperature and salinity vs number of training epochs.

The neural network with temperature overfits, as the validation loss is lower than the training loss, potentially due to too small a validation data set, which was reduced by using the k-fold validation method resulting in reduced test RMSE of 0.167 °C and 0.391PSU.

This is an introduction to the potential for the application of this process to satellite multispectral sensors, which unlike hyperspectral sensors, provide resolutions as high as 10 m (Sentinel-2) with global coverage, Caballero et al. (2020). This is much higher than that provided by typical salinity-measuring sensors, like microwave remote sensing, such as the European Space Agency's

Soil Moisture Ocean Salinity satellite (SMOS, 1 km spatial resolution), SMO (2020). The increase in resolution would enable the monitoring of complex water types, e.g. in coastal waters for differentiating river plumes, upwelling, evaporation, precipitation or vertical mixing. The decrease in image pixel size will also mean the algorithms would be able to monitor closer to the land sea boundary without the significant land adjacency effects, Rani et al. (2021), Muller-Karger (2018).

This study has demonstrated a relatively straightforward test case scenario due to the homogeneity of the water type and the excellent data quality provided by hyperspectral 1 nm band input data. The next steps involve applying the algorithm to more complex optical waters, characterized by greater variability in factors such as chlorophyll-a, coloured dissolved organic matter (CDOM), and suspended particulate matter (SPM).

The transition from hyperspectral data to multispectral satellite remote sensing data presents its own set of challenges, such as atmospheric effects and sensor-specific issues, which will require the use of sophisticated atmospheric correction techniques and more refined algorithms to ensure accurate and reliable predictions, White et al. (2024). While the development of NASA's Plankton, Aerosol, Cloud, ocean Ecosystem (PACE) mission, (NASA, 2025), marks an exciting stage in global ocean colour coverage (Groom et al., 2019), it also introduces some challenges. PACE provides global ocean colour coverage with a resolution of 1 km², which may not be ideal for monitoring the intricate dynamics of coastal estuaries. In contrast, sensors such as PRISMA (PRecursore IperSpettrale della Missione Applicativa), HICO (Hyperspectral Imager for the Coastal Ocean), and EMIT (Earth Surface Mineral Dust Source Investigation) are equipped with much higher spatial resolutions, making them more suitable for coastal environments where spatial detail is critical. These sensors have the ability to provide high-resolution, hyperspectral data at scales more relevant for precise monitoring of coastal ecosystems. As such this methodology could be applied across hyperspectral bands of different missions, producing high spatial and/ or temporal results when needed. By applying advanced algorithms to high-resolution satellite data, we can also validate satellite missions such as those from TRISHNA (Thermal infraRed Imaging Satellite for High-resolution Natural resource Assessment).

This paper highlights the importance of increasing the spatial resolution of satellite-derived sea surface temperatures (SSTs) and sea surface salinity (SSSs) in coastal regions, particularly in sensitive estuarine environments. Traditional global datasets often miss the fine-scale dynamics of coastal waters, but with higher resolution data, this approach enables more precise monitoring of coastal ecosystems, enhances climate models, and supports better marine resource management. Ultimately, this work contributes to improving our ability to track and respond to environmental changes in diverse marine environments.

CRediT authorship contribution statement

Solomon White: Writing – review & editing, Writing – original draft, Visualization, Validation, Software, Methodology, Investigation, Formal analysis, Data curation, Conceptualization.

Declaration of competing interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

Acknowledgements

Dr. Jose Luis Iriarte M, Universidad Austral de Chile for the salinity and temperature cruise data. Funded by the FONDAP-CONICYT, to the Chilean Government. Solomon White acknowledges Cefas as the PhD CASE partner. PhD in the SENSE CDT – funded by National Environmental Research Council (NERC), United Kingdom and UKSA, United Kingdom. ChatGPT AI was used to help structure the cover letter for submission to the journal.

Data availability

Data will be made available on request.

References

- Altieri, A.H., Gedan, K.B., 2015. Climate change and dead zones. *Global Change Biol.* 21 (4), 1395–1406.
- Ardeshir, N., Sanford, C., Hsu, D., 2021. Support vector machines and linear regression coincide with very high-dimensional features. *arXiv preprint arXiv: 2105.14084*.
- Behrenfeld, M.J., O'Malley, R.T., Siegel, D.A., McClain, C.R., Sarmiento, J.L., Feldman, G.C., Milligan, A.J., Falkowski, P.G., Letelier, R.M., Boss, E.S., 2006. Climate-driven trends in contemporary ocean productivity. *Nature* 444 (7120), 752–755.
- Binding, C.E., Bowers, D.G., 2003. Measuring the salinity of the Clyde Sea from remotely sensed ocean colour. *Estuar. Coast. Shelf Sci.* 57 (4), 605–611.
- Breitburg, D., Levin, L.A., Oschlies, A., Grégoire, M., Chavez, F.P., Conley, D.J., Garçon, V., Gilbert, D., Gutiérrez, D., Isensee, K., Jacinto, G.S., Limburg, K.E., Montes, I., Naqvi, S.W.A., Pitcher, G.C., Rabalais, N.N., Roman, M.R., Rose, K.A., Seibel, B.A., Telszewski, M., Yasuhara, M., Zhang, J., 2018. Declining oxygen in the global ocean and coastal waters. *Science* 359 (6371), eaam7240.
- Bricaud, A., Claustre, H., Ras, J., Oubelkheir, K., 2004. Natural variability of phytoplanktonic absorption in oceanic waters: Influence of the size structure of algal populations. *J. Geophys. Res.: Oceans* 109 (C11).
- Caballero, I., Fernández, R., Escalante, O.M., Mamán, L., Navarro, G., 2020. New capabilities of Sentinel-2A/B satellites combined with in situ data for monitoring small harmful algal blooms in complex coastal waters. *Sci. Rep.* 10 (1), 1–14.

- Cael, B.B., Chase, A., Boss, E., 2020. Information content of absorption spectra and implications for ocean color inversion. *Appl. Opt.* 59 (13), 3971.
- Casey, K.A., Rousseaux, C.S., Gregg, W.W., Boss, E., Chase, A.P., Craig, S.E., Mouw, C.B., Reynolds, R.A., Stramski, D., Ackleson, S.G., Bricaud, A., Schaeffer, B., Lewis, M.R., Maritorena, S., 2020. A global compilation of in situ aquatic high spectral resolution inherent and apparent optical property data for remote sensing applications. *Earth Syst. Sci. Data* 12, 1123–1139.
- Cover, T.M., Thomas, J.A., 2006. *Elements of Information Theory*, second ed. Wiley-Interscience, Hoboken, NJ.
- Darecki, M., Stramski, D., 2004. An evaluation of MODIS and SeaWiFS bio-optical algorithms in the Baltic Sea. *Remote Sens. Environ.* 89 (3), 326–350.
- Domingos, P., 2012. A few useful things to know about machine learning. *Commun. ACM* 55 (10), 79–88.
- European Space Agency, 2019. Sentinel-2 Products Specification Document. Technical Report, ESA.
- FONDECYT - National Fund for Scientific and Technological Development, 2023. FONDECYT Program.
- Fournier, S., Lee, T., Gierach, M.M., 2016. Seasonal and interannual variations of sea surface salinity associated with the Mississippi River plume observed by SMOS and Aquarius. *Remote Sens. Environ.* 180, 431–439.
- Giddy, I., Swart, S., du Plessis, M., Thompson, A.F., Nicholson, S.-A., 2021. Stirring of sea-ice meltwater enhances submesoscale fronts in the southern ocean. *J. Geophys. Res.: Oceans* 126, e2020JC016814, Received 22 SEP 2020, Accepted 2 MAR 2021.
- Groom, S.B., Sathyendranath, S., Ban, Y., Bernard, S., Brewin, B., Brotas, V., Brockmann, C., Chauhan, P., Choi, J.K., Chuprin, A., Ciavatta, S., Cipollini, P., Donlon, C., Franz, B.A., He, X., Hirata, T., Jackson, T., Kampel, M., Krasemann, H., Lavender, S.J., Pardo-Martinez, S., Melin, F., Platt, T., Santoleri, R., Skakala, J., Schaeffer, B., Smith, M., Steinmetz, F., Valente, A., Wang, M., 2019. Satellite ocean colour: Current status and future perspective. *Front. Mar. Sci.* 6 (JUL).
- Hsieh, 2009. Machine learning methods in the Environmental Sciences Neural Networks and Kernel.
- Hsu, T.-J., 2016. Sediment resuspension. In: Kennish, M.J. (Ed.), *Encyclopedia of Estuaries*. Springer Netherlands, Dordrecht, pp. 558–560.
- Hu, C., Montgomery, E.T., Schmitt, R.W., Muller-Karger, F.E., 2004. The dispersal of the Amazon and Orinoco River water in the tropical Atlantic and Caribbean Sea: Observation from space and S-PALACE floats. *Deep-Sea Res. Part II: Top. Stud. Ocean.* 51 (10–11 SPEC. ISS.), 1151–1171.
- Hughes, T.P., Anderson, K.D., Connolly, S.R., Heron, S.F., Kerry, J.T., Lough, J.M., Baird, A.H., Baum, J.K., Berumen, M.L., Bridge, T.C., 2018. Spatial and temporal patterns of mass bleaching of corals in the anthropocene. *Science* 359 (6371), 80–83.
- Ianson, D., Allen, S.E., Harris, S.L., Orians, K.J., Varela, D.E., Wong, C.S., 2003. The inorganic carbon system in the coastal upwelling region west of Vancouver island, Canada. *Deep-Sea Res. Part I: Ocean. Res. Pap.* 50 (8), 1023–1042.
- Jackson, L.C., Biastoch, A., Buckley, M.W., Desbruyères, D.G., Frajka-Williams, E., Moat, B., Robson, J., 2022. The evolution of the north atlantic meridional overturning circulation since 1980. *Nat. Rev. Earth & Environ.* 3 (4), 241–254.
- James, G., Witten, D., Hastie, T., Tibshirani, R., 2013. An introduction to statistical learning. In: *Encyclopedia of Machine Learning*. Springer.
- Li, G., Cheng, L., Zhu, J., et al., 2020. Increasing ocean stratification over the past half-century. *Nature Clim. Change* 10, 1116–1123.
- López, O.A.M., López, A.M., Crossa, J., 2022. Support vector machines and support vector regression. In: *Multivariate Statistical Machine Learning Methods for Genomic Prediction*. Springer.
- MacCreedy, P., Geyer, W.R., 2010. Advances in estuarine physics. *Annu. Rev. Mar. Sci.* 2 (1), 35–58.
- Minnett, P.J., Alvera-Azcárate, A., Chin, T.M., Corlett, G.K., Gentemann, C.L., Karagali, I., Li, X., Marsouin, A., Marullo, S., Maturi, E., Santoleri, R., Saux Picart, S., Steele, M., Vazquez-Cuervo, J., 2019. Half a century of satellite remote sensing of sea-surface temperature. *Remote Sens. Environ.* 233 (September).
- Morel, A., Prieur, L., 1977. Analysis of variation in ocean color. *Limnol. Oceanogr.* 22, 709–722.
- Muller-Karger, 2018. Satellite sensor requirements for monitoring essential biodiversity variables of coastal ecosystems. *Ecol. Appl.* 28 (3), 749–760.
- NASA, 2025. PACE Satellite Home. Accessed: 2024-06-18 URL <https://pace.gsfc.nasa.gov/>.
- Oliver, E.C., Donat, M.G., Burrows, M.T., Moore, P.J., Smale, D.A., Alexander, L.V., Benthuyzen, J.A., Bindoff, N.L., Hobday, A.J., Holbrook, N.J., 2018. Longer and more frequent marine heatwaves over the past century. *Nat. Commun.* 9 (1), 1324.
- O'Reilly, J.E., Werdell, P.J., 2019. Chlorophyll algorithms for ocean color sensors - OC4, OC5 & OC6. *Remote Sens. Environ.* 229, 32–47.
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M., Prettenhofer, P., Weiss, R., Dubourg, V., Vanderplas, J., Passos, A., Cournapeau, D., Brucher, M., Perrot, M., Duchesnay, E., 2011. Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.* 12, 2825–2830.
- Pollock, F.J., Lamb, J.B., Field, S.N., Heron, S.F., Schaffelke, B., Shedrawi, G., Bourne, D.G., Willis, B.L., 2014. Sediment and turbidity associated with offshore dredging increase coral disease prevalence on nearby reefs. *PLoS One* 9 (7).
- Rajaratnam, B., Roberts, S., Sparks, D., Dalal, O., 2016. Lasso regression: Estimation and shrinkage via the limit of gibbs sampling. *J. R. Stat. Soc. Ser. B Stat. Methodol.* 78 (1), 153–174.
- Rani, M., Masroor, M., Kumar, P., 2021. Remote sensing of Ocean and Coastal Environment – Overview. *Remote Sensing of Ocean and Coastal Environments*. Elsevier, pp. 1–15.
- Ranjam, J., Cook, J.A., 2018. LASSO regression. *Br. J. Surg.* 105 (10), 1348.
- Roettgers, R., Doerffer, R., Fischer, J., Stelzer, K., Reinart, A., Petersen, W., 2014. Temperature and salinity correction coefficients for light absorption by water constituents in seawater. *Opt. Express* 22 (21), 25093–25108.
- Rokem, A., Kay, K., 2020. Fractional ridge regression: a fast, interpretable reparameterization of ridge regression. *GigaScience* 9 (12), gaa133.
- Sen Gupta, A., Thomsen, M., Benthuyzen, J.A., Hobday, A.J., Oliver, E., Alexander, L.V., Burrows, M.T., Donat, M.G., Feng, M., Holbrook, N.J., Perkins-Kirkpatrick, S., Moore, P.J., Rodrigues, R.R., Scannell, H.A., Taschetto, A.S., Ummenhofer, C.C., Wernberg, T., Smale, D.A., 2020. Drivers and impacts of the most extreme marine heatwave events. *Sci. Rep.* 10 (1), 19359, Received 9 NOV 2020, Accepted 9 NOV 2020.
- Simis, S.G.H., Olsson, J., 2013. Unattended processing of shipborne hyperspectral reflectance measurements. *Remote Sens. Environ.* 135, 202–212.
- SMO, 2020. SMOS: The mission and the system. https://www.esa.int/Applications/Observing_the_Earth/SMOS (Accessed: 08/10/2020).
- Society, I.B., 1967. A Comparison of Some Methods of Cluster Analysis Author (s): J.C. Gower Published by : International Biometric Society Stable URL: <https://www.jstor.org/stable/2528417> REFERENCES Linked references are available on JSTOR for this article : reference, 23(4), 623–637.
- TriOS, 2024. TriOS sensors. Accessed: 2024-06-18 URL <https://www.trios.de/en/sensors.html>.
- Vali, A., Comai, S., Matteucci, M., 2020. Deep learning for land use and land cover classification based on hyperspectral and multispectral earth observation data: A review. *Remote Sens.* 12 (15).
- Wei, J., Wang, M., Ondrusek, M., Gilerson, A., Goes, J., Hu, C., Lee, Z., Voss, K.J., Ladner, S., Lance, V.P., Tuffillaro, N., 2023. Chapter 20 - Satellite ocean color validation. In: Nalli, N.R. (Ed.), *Field Measurements for Passive Environmental Remote Sensing*. Elsevier, pp. 351–374.
- White, S., Silva, T., Amoudry, L.O., Spyarakos, E., Martin, A., Medina-Lopez, E., 2024. The colours of the ocean: using multispectral satellite imagery to estimate sea surface temperature and salinity in global coastal areas, the gulf of Mexico and the UK. *Front. Environ. Sci.* 12.
- Wouthuyzen, S., Kusmanto, E., Fadli, M., Harsono, G., Salamena, G., Lekalette, J., Syahailatua, A., 2020. Ocean color as a proxy to predict sea surface salinity in the Banda Sea. *IOP Conf. Series: Earth Environ. Sci.* 618 (1).