

## Master Thesis

on the topic of

# Modelling and optimization of ship's fuel consumption using Random Forest Regression (RFR)

Submitted to the Faculty of Engineering  
of University Duisburg Essen

by

**Hibatul Wafi  
3021919**

Betreuer: M. T. Muhammad Fakhruriza Pradana  
1. Gutachter: Prof. Dr.-Ing. B. Noche  
2. Gutachter: Dr.-Ing. Alexander Goudz  
Studiengang: ISE General Mechanical Engineering  
Studiensemester: Summer semester 2023  
Datum: 04.05.2023

# Contents

<b>1</b>	<b>Introduction</b>	<b>5</b>
1.1	Thesis Objective . . . . .	6
1.2	Thesis Boundaries . . . . .	7
1.3	Thesis Contributions . . . . .	7
1.4	Thesis Structure . . . . .	7
<b>2</b>	<b>Theoretical Background</b>	<b>9</b>
2.1	Literature Review . . . . .	9
2.1.1	Modelling Approach for Ship Operation . . . . .	9
2.1.2	Overview of data source . . . . .	10
2.1.3	Review of ML approach to predict FOC . . . . .	11
2.1.4	Tree-Based Model as FOC model . . . . .	12
2.1.5	Conclusion of Literature Review . . . . .	14
2.2	Tree-based model . . . . .	14
2.2.1	Decision Tree . . . . .	14
2.2.1.1	Random Forest . . . . .	16
2.2.2	Extra-Trees (Extremely Randomised Trees) . . . . .	18
2.3	AIS Data . . . . .	18
2.3.1	Current Correction . . . . .	19
2.4	Weather data . . . . .	21
2.5	Calculation of Fuel Oil Consumption . . . . .	21
<b>3</b>	<b>Research Methodology</b>	<b>22</b>
3.1	Data Acquisition . . . . .	22
3.2	Data Preprocessing . . . . .	24
3.2.1	Data Cleaning . . . . .	24
3.2.2	Feature Selection . . . . .	27
3.2.3	Modelling . . . . .	30
3.2.4	Performance Metrics for Model Validation . . . . .	31
3.2.4.1	Coefficient of Determination ( $R^2$ ) . . . . .	32
3.2.4.2	Explained Variance (EV) . . . . .	32
3.2.4.3	Mean Absolute Error (MAE) . . . . .	32
3.2.4.4	Root Mean Square Error (RMSE) . . . . .	32
3.2.4.5	Median Absolute Deviation (MAD) . . . . .	33
3.2.5	Model Hyperparameter Optimisation . . . . .	33
3.2.5.1	Number of features . . . . .	34

3.2.5.2	Number of sample in a leaf node . . . . .	34
3.2.5.3	Depth of Tree . . . . .	34
3.2.5.4	Number of Trees . . . . .	35
3.2.6	Methodology Application . . . . .	35
3.2.7	Data Analysis . . . . .	37
3.2.8	Modelling . . . . .	40
3.2.9	Predicting STW . . . . .	40
<b>4</b>	<b>Result and Discussion</b>	<b>44</b>
4.0.1	Model Evaluation . . . . .	44
<b>5</b>	<b>Summary and Outlook</b>	<b>46</b>
<b>References</b>		<b>47</b>

# List of Tables

2.1	Structure of AIS data (IMO, 2015) . . . . .	20
3.1	Structure of fused dataset . . . . .	25
3.2	Structure of fused dataset . . . . .	30
4.1	Model performance . . . . .	44
4.2	Model performance . . . . .	45

# List of Figures

2.1 Example of partition space (Hastie et al., 2009) . . . . .	15
2.2 Example of partition tree (Hastie et al., 2009) . . . . .	15
2.3 Prediction of two Decision tree regression models (Géron, 2019) . . .	16
2.4 Regularising a Decision Tree regressor (Géron, 2019) . . . . .	17
3.1 Scheme of proposed methodology . . . . .	22
3.2 Shore based AIS Coverage based on data from AIS database Danish Maritime Authority (2023) . . . . .	26
3.3 Journey of the ship in a year . . . . .	27
3.4 Statistical distribution of wave heights Bretschneider (1965) . . . . .	28
3.5 Correlation Heat Map . . . . .	29
3.6 Visual illustration of k-folding, Grey shaded box represents the validation data while white box represents the training data . . . . .	31
3.7 Hyperparameter tuning of max_features . . . . .	34
3.8 Hyperparameter tuning of min_samples_leaf . . . . .	35
3.9 Hyperparameter tuning of max_depth . . . . .	35
3.10 Hyperparameter tuning of n_estimators . . . . .	36
3.11 Journey of the ship in June . . . . .	36
3.12 Histogram of the features . . . . .	38
3.13 Correlation Heat Map . . . . .	40
3.14 Correlation Heat Map . . . . .	43

# Chapter 1

## Introduction

The research on efficient ship operation is a direction that is being actively pursued by marine industry stakeholders. This research direction is motivated by the increasing price of fuel oil and stricter environmental regulations. Fuel onboard a ship is referred to as “bunkers” and it takes up a considerable portion of ship’s Operational Expenses (OPEX). It is known that bunker fuel takes up more than 50% of voyage costs and constitutes up to 75% of the ship’s total operating cost. It can be inferred that energy efficient ship operations that could reduce fuel consumption translate to an increase in profitability ([Stopford, 2009](#); [Ronen, 2011](#); [Bialystocki and Konovessis, 2016](#)). Additionally, efficient operation also means the reduction of Green House Gas (GHG) emissions. The most recent report by International Maritime Organisation indicated that GHG emissions from shipping make up 2.51% of global emissions ([IMO, 2020](#)). This alignment in motivation implies that through energy efficient ship operation, marine industry stakeholders gain economic benefits while adhering to stringent environmental regulations.

With that, maritime industry stakeholder actively searches for methods to ensure energy efficient operation. Two approaches are considered, namely technical solutions and operational solutions. Technical solutions involve modification to the vessel’s structure and power system. But these solutions are expensive, and it requires engineering innovations ([Yan et al., 2021](#); [Li et al., 2022](#)). Because of this, stakeholders look for cheaper solutions to achieve energy efficient operation. The answer for an inexpensive approach lies in optimisation of operational measures, it carries less cost, and it does not require initial investments. Several recommended solutions can be found in Ship Energy Efficiency Management Plan (SEEMP).

However, greater focus will be given in this thesis towards optimising ship speed as reduction of ship speed has the greatest impact on fuel consumption. Different studies indicated that fuel consumption is correlated through a third-order, non-linear function of the ship speed ([Wang and Meng, 2012](#); [Ronen, 2011](#); [Du et al., 2019](#)). The significant impact of ship speed on fuel consumption is further supplemented by reports and studies stating that reducing ship speed by about 2 – 3 knots could halve the operating cost of shipping companies ([Stopford, 2009](#); [Wijnolst et al., 2009](#)). For these reasons, slow steaming is the measure that is most widely adopted

by shipping operator.

While inexpensive, optimising operational measures is not an easy and trivial task. Several factors ranging from vessel operational performance to varying weather conditions make it challenging to model the ship speed. Some fuel consumption models, which are based on historical data and ship parameters, lack generalisation capabilities, and it is sensitive towards noisy data. To address this problem, recent research turns towards data-driven approach i.e. machine learning approach to predict ship speed and fuel consumption. These studies reported success in their modelling, citing good generalisation capability and low prediction errors. Despite these successes, maritime experts find it difficult to accept models based on data driven approach, as some data-driven models are complex as well as unintuitive and in some cases can violate basic physical knowledge of the vessel. The performance of the data-driven model is also greatly dependent on both data quantity and quality ([Yan et al., 2021](#); [Gkerekos et al., 2019](#)).

As such, prompted by volatility and ever-increasing bunker fuel price, developing a model that could accurately predict Fuel Oil consumption (FOC) could prove to be useful to maritime industry stakeholders. As stakeholders could make critical economical decisions at the most opportune moment without violating the stringent environmental regulations.

## 1.1 Thesis Objective

This thesis proposes an intuitive, data-driven modelling approach that considers varying ship state and environment conditions to predict fuel consumption. To ensure the abundance of data during modelling, this thesis utilise data fused between Automatic Identification System (AIS) and weather data.

To achieve this, Grey Box Model (GBM) approach is selected. Machine learning approach using random forest regressor (RFR) is considered to provide a certain degree of intuitiveness to predict ship over ground (SOG) over different journey periods using fused AIS and weather data. Predicted SOG is then converted to actual ship speed i.e. Speed Through Water (STW). STW will be used as the input for modelling of Fuel Oil Consumption (FOC), which is carried out through Holtrop-Mennen estimation method ([Holtrop and Mennen, 1978, 1982; Holtrop, 1984](#)), a power estimation method based on hydrodynamic laws which consider resistance forces exerted by environmental conditions.

The following Research Questions (RQs) could be raised during the development of the model :

- **RQ1:** What are the steps that should be taken to optimise the predictive performance of the model?

- **RQ2:** Is it feasible to fuse AIS data and meteorological data to accurately predict the ship's SOG and subsequently FOC of the ship?
- **RQ3:** Which approximations and empirical equations are suitable to estimate the resistance forces required to estimate the power required by the ship?

## 1.2 Thesis Boundaries

The following research boundaries are set throughout this thesis:

- Due to the continuous nature of the SOG, only the regression aspect of Random Forest (RF) will be considered.
- The focus of this work is a detailed study of the performance and possible optimisation configuration of different tree-based predictors for SOG. As such, an exhaustive comparison study between different types of machine learning models will not be performed.
- In the case study, the approximation for ship parameters and dimensions is based on a similar type of ship with nearly identical dimensions.

## 1.3 Thesis Contributions

The GBM approach using the fusion of AIS data and weather data provides the following contributions :

- Economical and independent data source.
- Robust modelling approach that requires minimal data pre-processing and minimal model configuration.
- Comprehensible model that adheres to physical principles and hydrodynamic laws of the vessel.

## 1.4 Thesis Structure

The thesis is organised with the following structure:

**Chapter 1** introduces the problem statement and described the objective and boundaries of the thesis. The novelty of this thesis is declared in this chapter.

**Chapter 2** The fundamental aspects of the methodologies used to develop the model will be explained in this chapter. Section 2.1 including literature review of relevant past and present research. The fundamentals of the tree-based model will be discussed in Section 2.2, basic explanation of the parameters used in AIS and weather data will be given in Section 2.3 and Section 2.4. Section 2.5 presents the empirical

formulas and parameters used to estimate fuel consumption used by the ship based on various literature studies.

**Chapter 3** discuss the methodology used to develop tree-based model used for SOG prediction. The discussion comprises analysis of training data, feature selection and reduction and selection of tuning parameters of the model. The methodology to estimate resistance for ship power estimation will be discussed in this chapter as well.

**Chapter 4**, the GBM model will be evaluated using appropriate performance metrics and their effectiveness will be discussed. The review of the strength and limitations concerning the GBM method will be discussed here.

**Chapter 5** The summary of this study and reflections on the research process will be presented here.

# Chapter 2

## Theoretical Background

### 2.1 Literature Review

The literature review in Section 2.1 presents past and present research on utilisation of machine learning methods to achieve energy efficient operation. The concept of different modelling approaches for ship operation will be discussed in Section 2.1.1. The generalisation performance of random forest in various research will be discussed in ???. Brief summary of the literature review is presented in Section 2.1.5.

#### 2.1.1 Modelling Approach for Ship Operation

According to [Haranen et al. \(2016\)](#) and [Coraddu et al. \(2017\)](#), the modelling strategies to predict fuel consumption are classified into three categories:

**White Box Models (WBM)** are based on *a priori* mechanistic knowledge and physical principles of the vessel's system. This means that the dimensions of the vessel's structure, design parameters, and propulsion plant configuration are known.

**Black Box Models (BBM)** are purely data driven, and it is developed using data from different sailing journey and historical observations. Contrary to WBM, this approach does not require detailed information on the vessel. This modelling approach can be further split into two categories. *Statistical Modelling* aims to find explanations for relationships between fuel consumption and different factors that affect it. *Machine Learning (ML) Modelling* focuses on the predictive capabilities of the model that could predict fuel consumption at different points in time.

**Grey Box Models (GBM)** fuse WBM and BBM into a single model that considers both *a priori* knowledge of the vessel and historical sailing data. This method aims to complement the performance of WBM and BBM.

Each of these strategies possesses its strength and limitations. WBMs are developed based on physical and hydrodynamics laws as well as theories of naval architecture,

it is transparent and comprehensible, making them the preferred model used by various shipping industries. However, the deterministic nature of WBMs causes them to have poor suitability and generalisability. This is mainly caused due to limited *a priori* knowledge of different vessel dimensions, parameters, and narrow application limits of principle dimensions and form parameters of the vessel. Subsequently, the inability of WBMs to add randomness makes it rigid and restrictive. ([Haranen et al., 2016](#); [Yan et al., 2021](#))

BBMs in general have a good fitting ability for training data and good predictive accuracy for unseen data. BBMs developed using machine learning approach can generalise better compared to BBMs that are based on statistical modelling ([Petersen et al., 2012a](#)). BBMs are purely data driven, which means BBMs do not require former knowledge of vessel principle dimensions and form parameters. With increasing amount of data, better generalisation performance and handling of noisy data should be expected in a BBM. However, for the same reason, the quality of BBM model is highly dependent on data quantity and quality. For BBMs based machine learning approach, the amount of data is a major factor in determining the effectiveness of machine learning ([Halevy et al., 2009](#)). Data driven approach means that BBMs neglect basic vessel physical knowledge and are generally complex making it challenging to analyse and explain. For these reasons, experts in shipping industries are critical of models that do not include basic vessel knowledge and those that violate concepts of the domain knowledge in serious ways ([Yan et al., 2021](#)).

Hence, GBMs are introduced to address the limitations of both WBMs and BBMs by combining the mechanistic knowledge of the ship and physical principles of the vessel's system with BBM models, which possess good predictive capability. Despite these advantages, [Yan et al. \(2021\)](#) noted that GBM approach is not a common approach, recent research to predict fuel consumption are mainly dominated by BBM approach, specifically BBM based on machine learning approach.

### 2.1.2 Overview of data source

The modelling of FOC using GBM requires both components of WBM and BBM. For the BBM modelling part using machine learning approach, it is especially important to ensure sufficient amount of good quality data to be available for model training to ensure precise and accurate training of the model ([Halevy et al., 2009](#)).

It can be summarised, that the modelling of FOC use the following types of data source ([Yan et al., 2021](#)):

**(Daily) Noon Report:** Daily reports manually filed by ship's chief engineer and sent by the ship's masters to the shipping company and shore management. The reports include informations on types of daily fuel consumption, basic voyage information (e.g. ship location, load condition), sailing behaviour information e.g. (average sailing speed, average engine revolution per minute (RPM)), as well as sea and weather

conditions. While it provides relevant information regarding the ship operation, the inherent problem of daily and manual data entry means that the quality and quantity of data cannot be guaranteed.

**Sensor Data:** Data obtained from installed sensor onboard the vessel. This may include fuel flow sensors, Global Positioning System (GPS) receiver and wind speed sensors are among the possible sensors that can be installed onboard a vessel. Sensor data address the issues of data quantity from noon report, as pointed out in the study by [Gkerekos et al. \(2019\)](#) for the prediction of daily FOC. The machine learning models, which are produced by the Automated Data Logging and Monitoring (ADLM) system outperforms the models that used noon data for their training by 5 – 7% for a collection period of 3 months of the ADLM system and 2.5 years for the noon data. However, installing onboard sensors may be complex and costly ([Jóan Petur Petersen, 2011](#)) and the resulting sensor data will need to be handled properly to account for error in the sensors.

**AIS Data:** Apart from its intended use as collision avoidance system AIS data have seen potential usage in the field of scientific research. In the third Green House Gas (GHG) study by Smith et al. ([Smith et al., 2015](#)), uses AIS to estimate global shipping emission inventories. Rakke ([Rakke, 2016](#)) proposed a methodology termed ECAIS to calculate ship emissions based on the fuel consumption from AIS data. Through Holtrop-Mennen approximation and literature approximation, the ship's power propulsion can be determined which is subsequently used to predict specific fuel consumption. Wen et al. ([Wen et al., 2017](#)) attempted to minimise the Energy Efficiency Operational Indicator (EEOI) using green routing. Recent research by Kim et al. ([Kim et al., 2020](#)) used publicly accessible AIS data, ship static data and environmental data to estimate EEOI without requiring the actual FOC. The study used big data technology as public data are of large capacity. Generally, the study using AIS data is done to achieve independence from the need to use commercial database. The detail of AIS data will be discussed in Section 2.3

### 2.1.3 Review of ML approach to predict FOC

Modelling of FOC using *machine learning* approach generally focus on prediction of unseen data. The general framework usually include collection and preprocessing of ship operational data, training and validation of the model, and evaluation and selection of the most appropriate model. Some machine learning models allow further hyperparameter tuning of the model and in case of data rich environment, the data can be further split into test data to further validate the performance of machine learning model.

The study by [Yan et al. \(2021\)](#) indicated that the majority of recent research that uses machine learning approach employ ANN as the model to predict FOC. ANN models are powerful models capable of modelling nonlinear data which are based on theories on how the brain works. The outcome is modelled by intermediate set

of unobserved variables known as hidden layer. ([Kuhn and Johnson, 2013](#)). Back propagation neural networks, Multi Level Perceptron (MLP), and wavelet neural networks are some examples of ANN model subclasses.

ANN has shown respectable performance in its attempt to predict FOC. [Petersen et al. \(2012b\)](#) reported Root Mean Square Error of 47.2 L/h for the fuel flow i.e. FOC. To put this into context, the fuel flow in their case study fluctuates between 1000–2500 L/h. [Bal Beşikçi et al. \(2016\)](#) considered sailing speed, trim, wind, sea effects, propeller pitch, and engine rotation speed as input variables to predict FOC per hour and achieved model fit score of  $R^2 = 0.759$  in test set. Other studies also reported similar range of results using ANNs ([Yan et al., 2021](#)).

However, the development of ANN models is a challenging task. ANN models tend to overfit when there is shortage of data, as such, regularisation is necessary to improve model performance. The balancing process during regularisation is a demanding task and unsuitable regularisation may lead to counterintuitive prediction results. Adding layers is computationally expensive, and it does not always guarantee promising results ([Hastie et al., 2009](#)). Additionally, in machine learning terms, ANN is classified as a black box model, which makes it unintuitive and lacking in interpretability ([Géron, 2019](#)), this particular limitation cause shipping industry expert generally reluctant to accept the model generated using machine learning approach.

#### 2.1.4 Tree-Based Model as FOC model

Concerning interpretability, modelling approaches such Linear Regression (LR), KNN and tree-based models have shown superior interpretability in comparison to ANNs. LR can explain the effect of each input variable on the output through the coefficients. KNN searches for the nearest neighbour and their closeness is evaluated through distance measurement algorithms such as Euclidean distance. Additionally, LRs and KNNs also offer easy implementation and adequate explainability. However, both approaches suffer from sensitivity to outliers and noise in data

This brings us to tree-based model, a powerful and intuitive modelling approach capable of performing classification tasks for discrete data and regression tasks for continuous data. The main idea behind tree-based model follows the principle of partition space, the data points are split into their respective segment according to a certain threshold. The split is performed until a certain stopping rule is applied or when there are no more data points available for splitting and this splitting process for a single tree can be visualised using binary tree representation. ([Hastie et al., 2009](#)).

However, a single decision tree is prone to overfitting and sensitive to outliers, but this limitation can be resolved through regularisation or by creating an ensemble of independent decision trees, known as Random Forest (RF). In Random Forest Re-

gressor (RFR), the prediction is the average of the prediction across the decision trees **Breiman (2001)**. The detailed principle of tree-based regressor will be discussed in greater detail in Section 2.2.

There have been various research that considered tree based model to model FOC:

**Soner et al. (2018)** used the ferry dataset from **Petersen et al. (2012b)** to predict FOC using tree-based model, which includes bagging, random forest (RF), and bootstrap. From the test dataset, the random forest model achieved RMSE of 43.5 L/h for the fuel consumption. Which suggested improvement from ANN model from the study of **Petersen et al. (2012b)**.

**Yan et al. (2020)** used random forest (RF) model to predict FOC for a voyage of a dry bulk ship using ship operational data i.e. ship noon data and sea and weather data from noon report and EMCWF. The prediction model considered ship sailing speed, total cargo weight and meteorological conditions and RF model obtained mean absolute percentage error (MAPE) of 7.91% for the FOC. The RF model displayed superior result in comparison to Decision Tree Regressor (DTR), ANN, LASSO, and SVR.

The advantage of tree-based model is further highlighted by **Gkerekos et al. (2019)**. The study compared the performance of different machine learning models to predict ship's FOC per day using both noon data and automated data logging and monitoring (ADLM) system from a bulk carrier. This research concludes that tree-based model displayed good prediction performance on both noon data and sensor-based data. ETR achieved remarkable model fit score of 89% using the noon data and 97% when using the data from ADLM system, outperforming ANN, SVR, and RFR models.

**Li et al. (2022)** performed more extensive research on the effects of data fusions between meteorological data, ship voyage data, and AIS data on different machine learning models to predict the ship's FOC. The study classified ETR and RFR as tree-based model which is produced by *bagging ensemble strategy*. While AdaBoost (AB), Gradient Tree Boosting (GB), XGBoost(XG) and LightGBM (LB) are classified as tree-based models produced by *boosting ensemble strategy*. The study recommends all tree-based models that are produced by *boosting ensemble strategy* and ETR to be used to model energy efficient operation. Additionally, RFR shows the best robustness among the proposed model in the study.

**Abebe et al. (2020)** attempted to use machine learning approach to predict SOG of the ship. In this study, AIS data and noon-report weather data from 14 tracks and 62 ships are used for model training. The model considered the ship draught, ship dynamic information, tonnage, and environmental conditions. The result of this study exhibited the feasibility of using AIS data and meteorological data to predict SOG of the ship. The results also further indicated the strength of tree-based model, on test dataset, ETR achieved the best result with model fit of 98.47% and RMSE of 0,234 knots. It is also reported that ETR achieved better performance with about half of

the computational cost of RFR.

### 2.1.5 Conclusion of Literature Review

This literature review described the capability of Random Forest Regressor to predict fuel consumptions and ship speed, irrespective of data source and type of data used. Promising results from different performance measures across different literatures indicated the capability random forest model as predictor. As such, this thesis aims to find optimisation possibilities to extract maximum prediction performance from random forest. Due to the nonlinear, third order function estimate of fuel consumption (Ronen, 1982, 2011). Accurate prediction of ship speed is paramount to ensure optimal ship operation resulting in increase of profitability.

## 2.2 Tree-based model

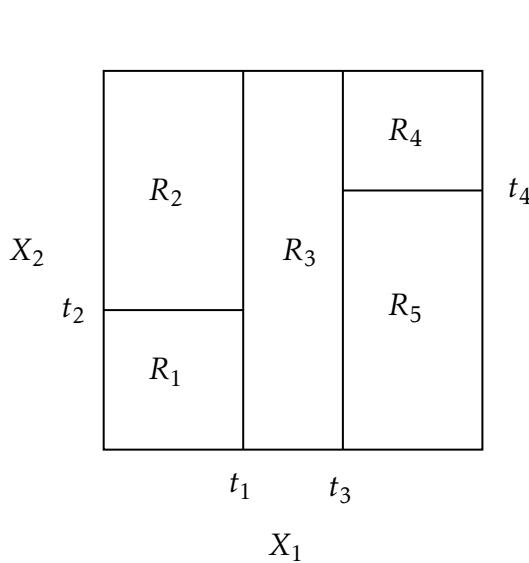
Random forest belongs to the family of tree-based model and its functional principle stems from decision tree. Decision tree is a non-parametric model that can perform both classification and regression tasks for discrete variable and continuous variable. It is a powerful algorithm, capable of fitting complex datasets. Tree-based model requires very little to no data pre-processing (Géron, 2019; Hastie et al., 2009). To grasp the concept of random forest, The principle working of decision tree will be introduced in depth in Section 2.2.1. It is then followed by Section 2.2.1.1 which presents the principle function behind random forest. Brief explanation for Extra-Trees, method introduced to further improvise random forest, will be presented in Section 2.2.2.

### 2.2.1 Decision Tree

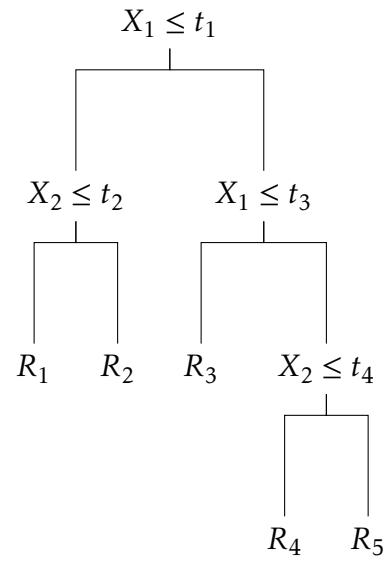
Decision tree is a white box model<sup>1</sup> (Géron, 2019). In machine learning sense, this means that the model is intuitive, and the structure of the model is interpretable. Thus, the structure of the model can be analysed in detail. To train Decision Trees, Scikit-Learn (Fabian Pedregosa et al., 2011) uses the *Classification and Regression Tree* (CART) algorithm (Breiman et al., 2017). Partition space shown by Figure 2.1 are used to illustrate the decision of CART algorithm. This process can be alternatively represented by the binary tree of Figure 2.2, observation that satisfies the condition are assigned to the left branch and the opposite is assigned to the right branch. The binary tree representation can be especially helpful when multiple input variables are involved, as the responses can be represented by a single tree (Hastie et al., 2009).

---

<sup>1</sup>This is not to be interchanged with the definition described by Haranen et al. (Haranen et al., 2016) regarding modelling of ship operation.



**Figure 2.1:** Example of partition space  
(Hastie et al., 2009)



**Figure 2.2:** Example of partition tree  
(Hastie et al., 2009)

Now, we need to understand the principle of selection for the feature  $k_t$  and threshold  $t_k$ . We shall first start with the principle of selection of the threshold  $t_k$ ; Assuming a case with single feature  $k$  and response  $y$ , with  $m$  data points. The algorithm starts by looking for possible thresholds. This is determined by calculating the splitting value.<sup>2</sup>. Then, the mean of the response  $y$  of partition space  $S_1$  and  $S_2$  is calculated as seen in Figure 2.3.

This step is then followed by calculating the sum of squared error (SSE) of each data points in partition space  $S_1$  and  $S_2$  and dividing it by the number of data points  $m_{S_1}$  and  $m_{S_2}$  respectively to obtain the MSE. Subsequently, the MSE from the respective partition space  $S_1$  and  $S_2$  is summed. The process is then recursively repeated until a threshold  $t_k$  that produce minimum sum of MSE is determined. This algorithm is defined by the following cost function  $J(k, t_k)$ , with  $\hat{y}_{S_i}$ , being the mean of the response,  $y_{S_i}$ , in partition space  $S_i$ . (Géron, 2019; Kuhn and Johnson, 2013):

$$\text{MSE}_{S_i} = \frac{1}{m_{S_i}} \text{SSE}_{S_i} \quad \text{where } i = (1, 2) \quad (2.2.1)$$

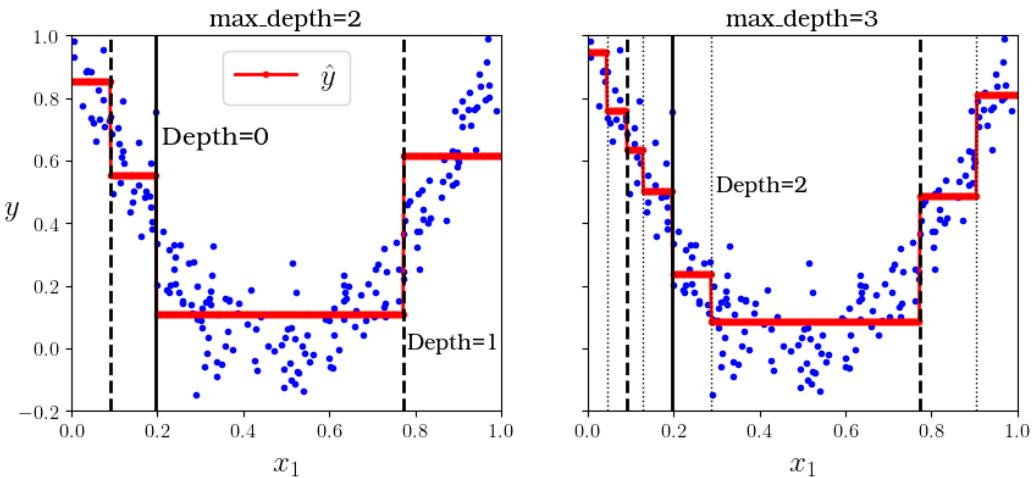
$$J(k, t_k) = \frac{1}{m_{S_1}} \text{SSE}_{S_1} + \frac{1}{m_{S_2}} \text{SSE}_{S_2} \left\{ \begin{array}{l} \text{SSE}_{S_i} = \sum_{i \in S_i} (\hat{y}_{S_i} - y_{S_i})^2 \\ \hat{y}_{S_i} = \frac{1}{m_{S_i}} \sum_{i \in S_i} y \end{array} \right. \quad (2.2.2)$$

Once complete, then the partition space is further split into two more regions and this process is recursively continued until a stopping rule is applied. The stopping rule are either when the tree reaches the maximum depth, (This is controlled by the parameter `max_depth` in Scikit-Learn), or when it cannot find a split that can

<sup>2</sup>For example, suppose there are data points at  $k = [0.2, 0.4]$ , then the splitting value is the value in between, i.e.,  $t_k = 0.3$

further reduce MSE. This best split also corresponds to the best possible fit to the predicted value.

Same principle is also applied when multiple features are present. Consider there are  $k_t$  features, then for each respective features  $k_1, k_2, \dots, k_t$ , The MSE for each of the features is calculated using the cost function  $J(k, t_k)$ . The feature that can *minimise* the cost function will be selected as the root of the tree. The tree is then grown further by recursively repeating this process (Hastie et al., 2009; Géron, 2019).

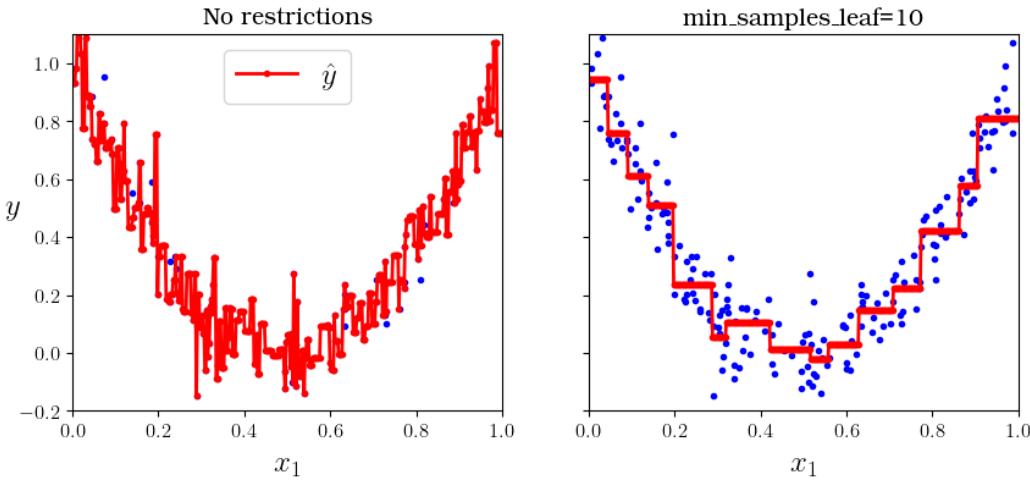


**Figure 2.3:** Prediction of two Decision tree regression models (Géron, 2019)

While powerful, decision tree suffers from overfitting when the model is unconstrained. Decision tree makes very few assumptions regarding the training data. Therefore, it will adapt to the training data and fitting it very closely (Géron, 2019). Additionally, an individual tree tends to be unstable, when the data is altered, a completely different set of splits might be found(Hastie et al., 2009; Kuhn and Johnson, 2013). Therefore, it is necessary to regularise i.e., restrict the decision tree's freedom during the training. Overfitting could be reduced by controlling how deep the tree can grow through the `max_depth` parameter. Additionally, setting the amount of minimum number of samples a leaf node has, through `min_samples_leaf` can alleviate overfitting as well, as shown in Figure 2.4. However, to address the fundamental drawbacks of decision tree, we shall look into random forest.

### 2.2.1.1 Random Forest

To understand random forest, the concept of ensemble method shall first be understood. Ensemble is defined as group of predictors such as classifier or regressor. Predictions are aggregated across multiple predictors, for regression task, the prediction is the average across the predictors. This principle is applied to random forest, a group of decision trees is trained on different random set of training data. For regression task, this means the prediction value is the average of the prediction across the decision trees. Such ensemble of decision trees is called **Random Forest** (Hastie



**Figure 2.4:** Regularising a Decision Tree regressor ([Géron, 2019](#))

[et al., 2009; Breiman, 2001; Tin Kam Ho, 1995](#)).

Ensemble methods achieve the best performance when the predictors are as independent to one another. In statistical sense, this can be achieved by reducing correlation among the trees. This can be realised by adding randomness during tree construction process. For this purpose, random forest utilises *bagging* ([Breiman, 1996](#)) method (short for *bootstrap aggregating*) during the training process. First, bootstrap sample is created, this means that a sample of the dataset is randomly selected and allowed to appear more than once. This sampling technique is referred to as sampling *with replacement*. Once predictors are trained, then the prediction of the new instance is *aggregated* across the predictors. ([Kuhn and Johnson, 2013; Hastie et al., 2009; Géron, 2019](#))

To add further randomness, random forest involves random selection of input features  $k$  that are considered to split the tree. This means that the feature  $k$  that will be used to split the tree is selected from this random subset of feature. The selection for the best feature to be used as the root of the tree and its subsequent node, as well as the stopping rule for the tree's growth is similar to that of decision tree. ([Kuhn and Johnson, 2013; Hastie et al., 2009; Géron, 2019](#))

These measures introduced in random forest address the tendency of decision tree to overfit. In fact, the instability of decision tree mentioned in Section 2.2.1 is exploited in random forest to gain randomness during construction of the tree. Experience from Hastie et al. ([Hastie et al., 2009](#)) shown that random forest requires minimal parameter tuning to achieve satisfactory performance while Kuhn et al. ([Kuhn and Johnson, 2013](#)) reported that tuning parameter does not have a drastic effect on performance.

However, what random forest gains in predictive performance, loses in interpretability.

ity. Random forest is considered as Black Box Model (BBM) ([Géron, 2019](#)).<sup>3</sup> The randomness means that it is challenging to analyse and describe the decisions made during the selection of the samples and during the selection of the input features. Nevertheless, the interpretability of a single tree in a random forest still holds. As it is still possible to traverse through the tree to reach the predicted value.

### 2.2.2 Extra-Trees (Extremely Randomised Trees)

Additionally, extra-trees (Extremely Randomised Trees) is introduced by Geurts et al. ([Geurts et al., 2006](#)) to further randomise random forest. The key difference lies on how each split is selected; in extra-trees each tree split is selected in random instead of searching for the best split. This technique saves computational power, as searching for best split is one of the tasks that takes up most computational power ([Géron, 2019](#)). Extra-trees also do not bootstrap the samples, which mean it samples *without replacement*.

## 2.3 AIS Data

Automatic Identification System (AIS) is an automated tracking system onboard ships to automatically transmit information about the ship to other ships and coastal authorities. As part of the revised new chapter V of SOLAS<sup>4</sup> regulation. In 2000, International Maritime Organization (IMO) requires installation of AIS class A equipment on all ships of 300 gross tonnage and upward engaged on international voyages, cargo ships of 500 gross tonnages and upwards not engaged on international voyages and all passenger ships irrespective of size. This requirement is then made compulsory to all ships by 2004. ([Rakke, 2016](#); [IMO, 2015](#))

AIS uses Very High Frequency (VHF) with special protocol for communication system for information exchange between the ships. This information will be received by either ships directly, buoys, land based station and satellites. The information transmitted by AIS is distinguished into three different types. **Static information** which is entered into the AIS on installation. **Dynamic information**, which is automatically updated from the ship's sensors connected to AIS and **voyage-related information**, which might need to be manually entered and updated during the voyage. The structure of the AIS data that is relevant to this thesis is summarised in Table 2.1([IMO, 2015](#)).

AIS is also further differentiated by its equipment class. The classification is based on the reporting interval and the type of information that is conveyed. **Class A** autonomously report their position within 2-10 seconds interval, depending on the

---

<sup>3</sup>Again, not to be interchanged with the definition described by Haranen et al. ([Haranen et al., 2016](#)) regarding modelling of ship operation.

<sup>4</sup>International convention for the Safety of Lives at Sea

state of ship's movement. The reporting interval is less frequent at 3 minutes, When the ship is at anchor or moored and moving slower than 3 knots. Class A AIS is also capable of sending safety related information, meteorological and hydrological data, electronic broadcast to mariners and marine safety messages. **Class B** reports at longer interval and at a lower power. They can only receive safety related messages, not send them. (Rakke, 2016; IMO, 2015)

### 2.3.1 Current Correction

As indicated in Table 2.1, the speed shown in AIS is the speed over ground (SOG). However, for calculation of ship's fuel consumption, the actual speed i.e. speed through water (STW) is required. This can be achieved by correcting the SOG for the current speed, in consideration of the research by Zhou et al. (Zhou et al., 2017) which shows the impact of current on ship's SOG. This correction is performed by considering the current speed  $V_c$  and the direction of the current  $\gamma$  *with respect to True North*. In principle, STW will be greater than SOG, when the current is moving against the current as the ship tries to compensate for the current to maintain the SOG. Similarly, the STW will be greater than the SOG when the current is moving in the same direction of the ship movement.

To calculate the correction, this study will adopt the methodology proposed by Kim et al. (Kim et al., 2020) and Yang et al. (Yang et al., 2020). The  $x$  and  $y$  component of SOG can be obtained through vector decomposition using the ship's heading angle  $\alpha$  *with respect to True North*. Similar vector decomposition is also performed for current speed  $V_{\text{current}}$ , it is resolved with current direction  $\gamma$  *with respect to True North*:

$$V_{\text{SOG}}^x = V_{\text{SOG}} \cdot \sin(\alpha) \quad (2.3.1)$$

$$V_{\text{SOG}}^y = V_{\text{SOG}} \cdot \cos(\alpha) \quad (2.3.2)$$

$$V_{\text{current}}^x = V_{\text{current}} \cdot \sin(\gamma) \quad (2.3.3)$$

$$V_{\text{current}}^y = V_{\text{current}} \cdot \cos(\gamma) \quad (2.3.4)$$

Then the resulting equation to determine STW, including the current compensation, is given by:

$$V_{\text{STW}}^x = V_{\text{SOG}}^x - V_{\text{current}}^x \quad (2.3.5)$$

$$V_{\text{STW}}^y = V_{\text{SOG}}^y - V_{\text{current}}^y \quad (2.3.6)$$

$$V_{\text{STW}} = \sqrt{(V_{\text{STW}}^x)^2 + (V_{\text{STW}}^y)^2} \quad (2.3.7)$$

Information Item	Description
<b>Static</b>	
MMSI	MMSI number of vessel
Callsign	Callsign of vessel
Name	Name of the vessel
IMO	IMO number of the vessel
Length	Length of vessel
Width	Width of vessel
Ship Type	Describes the AIS ship type of this vessel
<b>Dynamic</b>	
Ship's position	Automatically updated from position sensor connected to AIS. Longitude and Latitude.
Position time stamp in UTC	Automatically updated from ship's main position sensor. Format: DD/MM/YYYY HH:MM:SS
Course over Ground (COG)	<i>If available</i> , automatically updated from ship's main position sensor connected to AIS.
Speed Over Ground (SOG)	<i>If available</i> , automatically updated from the position sensor connected to AIS.
Heading	Automatically updated from the ship's heading sensor connected to AIS
Navigational status	Navigational status information has to be manually entered by the Officer on Watch (OOW) and changed as necessary. For example : “underway by engines”, “engaged in fishing”, “at anchor”.
Rate of Turn (ROT)	<i>If available</i> , Automatically updated from the ship's ROT sensor or derived from the gyro.
<b>Voyage Related</b>	
Ship's draught	To be manually entered at the start of the voyage using the maximum draft for the voyage and amended as required
(Hazardous) Cargo Type	Type of cargo from AIS message.
Destination and ETA	To be manually entered at the start of the voyage and kept up to date as necessary.

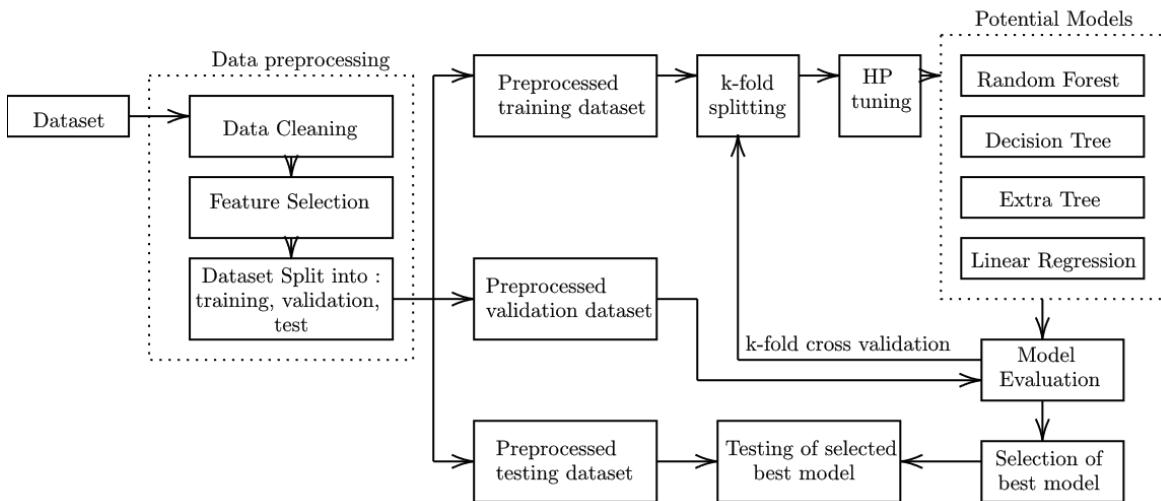
Table 2.1: Structure of AIS data (IMO, 2015)

**2.4 Weather data****2.5 Calculation of Fuel Oil Consumption**

# Chapter 3

## Research Methodology

In this chapter the methodology used to develop random forest model will be discussed. The details of fusion between AIS data, ECMWF and CMEMS data source used for training the model will be presented in Section 3.1. Suitable methodology application during data pre-processing will be described in Section 3.2. The selection for appropriate, domain knowledge based, feature selection will be explained in Section 3.2.2. The selection of the most optimal model hyperparameter for different tree-based model will be explained in Section 3.2.5. Different performance metrics is used to validate the model's generalisation capability, The underlying principle of the metrics is elaborated in Section 3.2.4. Summary of methodology application in this study is summarised in Section 3.2.6 and visually represented in figure Figure 3.1.



**Figure 3.1:** Scheme of proposed methodology

### 3.1 Data Acquisition

For the purpose of model training, 2021 AIS data from Ro-Ro ferry ship Hammershus is collected. The shore-based AIS data is made available by Danish Maritime

Authority which tracked her journey between ports of Køge, Rønne, Ystad and Sassnitz and structured according to Table 2.1. The AIS data is fused with weather data from ECMWF<sup>1</sup> with temporal resolution of 1 hour at granularity of 0.25° (longitude) x 0.25° (latitude), data from ECMWF provides information for wind, waves and seawater temperature. The information for current is obtained from CMEMS<sup>2</sup> with temporal resolution of 3 hours at granularity of 0.25° (longitude) x 0.25° (latitude).

The resulting fusion resulted in dataset with temporal resolution of 1 hour. Some information static information from the AIS data which only indicated the ship's identity are excluded. This includes ship's MMSI, Callsign, Name, IMO and Navigational Status. Additionally, information of the ship's Rate of Turn (ROT) is not available in this case. The weather information is synchronised so that the wind, waves, seawater temperature and sea current belongs to the same weather grid with same temporal resolutions.

The features (1) wind direction, (2) swell direction, (3) and wind wave direction are oriented to true north. However, to reflect the actual direction of weather effects that are acting on the ship, these features are converted to true direction; where true direction is defined as the direction of weather effect with respect to the bow of the ship. The value ranges between 0° and 180°. Subsequently, through vector decomposition, the northward and eastward wind velocity is converted to absolute wind speed and wind direction *with respect to True North*,  $\varphi$ :

$$V_{\text{wind}} = \sqrt{(V_{\text{wind}}^N)^2 + (V_{\text{wind}}^E)^2} \quad (3.1.1)$$

$$\varphi = \begin{cases} 360 - \arctan\left(\frac{V_{\text{wind}}^E}{V_{\text{wind}}^N}\right) & \text{if } V_{\text{wind}}^E > 0 \wedge V_{\text{wind}}^N < 0 \\ 180 - \arctan\left(\frac{V_{\text{wind}}^E}{V_{\text{wind}}^N}\right) & \text{if } V_{\text{wind}}^E < 0 \wedge V_{\text{wind}}^N > 0 \\ 270 - \arctan\left(\frac{V_{\text{wind}}^E}{V_{\text{wind}}^N}\right) & \text{if } V_{\text{wind}}^E > 0 \wedge V_{\text{wind}}^N > 0 \\ \arctan\left(\frac{V_{\text{wind}}^E}{V_{\text{wind}}^N}\right) & \text{otherwise} \end{cases} \quad (3.1.2)$$

Similarly, information of Northward and Eastward current Velocity is converted to absolute current speed and current direction *with respect to True North*  $\gamma$ .

$$V_{\text{current}} = \sqrt{(V_{\text{current}}^N)^2 + (V_{\text{current}}^E)^2} \quad (3.1.3)$$

---

<sup>1</sup>European Centre for Medium-Range Weather Forecast

<sup>2</sup>Copernicus Marine Environment Monitoring Service

$$\gamma = \begin{cases} 360 - \arctan\left(\frac{V_{\text{current}}^E}{V_{\text{current}}^N}\right) & \text{if } V_{\text{current}}^E < 0 \wedge V_{\text{current}}^N > 0 \\ 180 - \arctan\left(\frac{V_{\text{current}}^E}{V_{\text{current}}^N}\right) & \text{if } V_{\text{current}}^E > 0 \wedge V_{\text{current}}^N < 0 \\ 270 - \arctan\left(\frac{V_{\text{current}}^E}{V_{\text{current}}^N}\right) & \text{if } V_{\text{current}}^E < 0 \wedge V_{\text{current}}^N < 0 \\ \arctan\left(\frac{V_{\text{current}}^E}{V_{\text{current}}^N}\right) & \text{otherwise} \end{cases} \quad (3.1.4)$$

This conversion is performed as the information of current speed and current direction,  $\gamma$ , is necessary to perform the correction formula shown in Equation (2.3.5) and Equation (2.3.6). However, for training purpose, this feature will not be considered. Instead, the true current direction and true wind direction will be considered. The initial structure have 27 features, 9 AIS features and 18 weather features. The structure of the initial dataset i.e. before data preprocessing and feature selection, is summarised in Table 3.1

## 3.2 Data Preprocessing

This section presents the steps taken to during data preprocessing. The dataset will be first subjected to data cleaning which include identification of anomalies and missing values, the steps are explained in Section 3.2.1. Boundary condition is then applied to ensure that the model represent operating condition at steady state. Using domain knowledge, appropriate features are selected and discarded to ensure the model obeys shipping domain knowledge. This dataset is to be split into training, validation and test dataset. These steps will be further elaborated in Section 3.2.2.

### 3.2.1 Data Cleaning

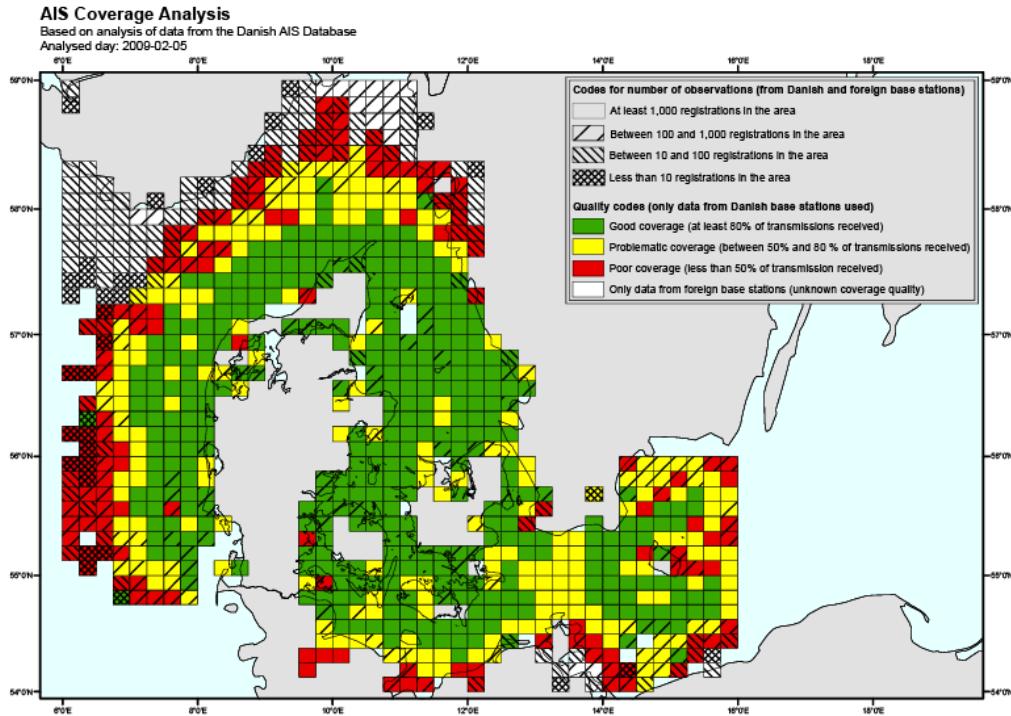
The journey between the port of Køge, Rønne, Ystad and Sassnitz is plotted using QGIS<sup>3</sup>. The plot of the journey is shown in Figure 3.3, it can be seen, that the journey between Rønne and Sassnitz is not represented completely. As in this information is missing due to poor coverage in the area between Sassnitz and Rønne. This is shown by the plot shown in Figure 3.2. Therefore, the data plot for the journey between Sassnitz and Rønne will be excluded. Basic threshold of decimal degrees of 55.04° N for latitude is applied, this threshold will exclude the journey between Sassnitz and Rønne.

In its initial state, the dataset contains 7453 data points which described the journey of the ship in one year. The initial data points represented all navigational status of the ship, which include “mooring”, “anchoring” and “underway using engine”. This is clearly observed in the histogram for the SOG distribution in figure BLALA. To ensure that the dataset represents the actual operating condition of ship in steady state, a threshold of 5 knots is applied. SOG can vary due to changing sea state,

<sup>3</sup><https://qgis.org/en/site/>, QGIS is a free and open source geographic information system

Feature	Feature Name
<b>AIS data</b>	
Position Time Stamp [DD/MM/YYYY HH:MM:SS]	Time
Latitude [ $^{\circ}$ ]	LAT
Longitude [ $^{\circ}$ ]	LON
Width [m]	width
Length [m]	length
SOG [Knots]	sog
COG [m/s]	cog
Heading [ $^{\circ}$ ]	heading
Draught [m]	draught
<b>Weather Data (0.5° Granularity)</b>	
Wind Speed [m/s]	windspeed
True North Wind Direction, $\varphi$ [ $^{\circ}$ ]	truenorthcurrentdir
Air Temperature Above Oceans [K]	oceantemperature
Maximum Wave Height [m]	waveheight
Swell Period [s]	swellperiod
Wind Wave Period [s]	windwaveperiod
Wave Period [s]	waveperiod
Sea Surface Temperature [K]	surftemp
Combined Wind Wave Swell Height [m]	windwaveswellheight
Swell Height [m]	swellheight
Wind Wave Height [m]	windwaveheight
Current Speed [m/s]	curspeed
True North Current Direction $\gamma$ [ $^{\circ}$ ]	truenorthcurrentdir
True Wind Direction [ $^{\circ}$ ]	truewinddir
True Current Direction [ $^{\circ}$ ]	truecurrentdir
True Swell Direction [ $^{\circ}$ ]	trueswelldir
True Wind Wave Direction [ $^{\circ}$ ]	truewindwavedir
True Wave Direction [ $^{\circ}$ ]	truewavedir

Table 3.1: Structure of fused dataset



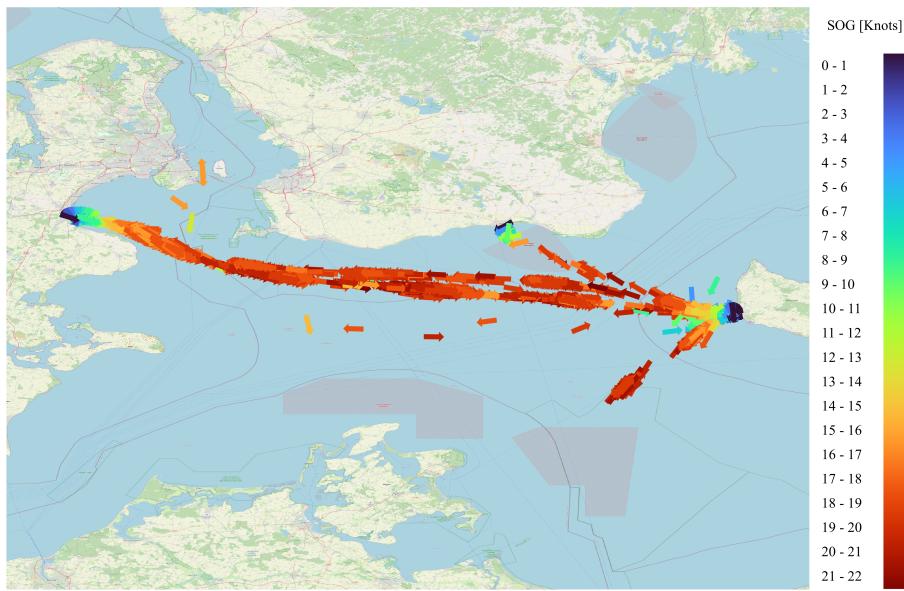
**Figure 3.2:** Shore based AIS Coverage based on data from AIS database [Danish Maritime Authority \(2023\)](#)

but it can also be reduced by the ship's operator around the port when it departs from port of origin or arriving at port of arrival. Any data points with SOG less than 5 knots will be discarded which is considered as manoeuvring [Abebe et al. \(2020\)](#). After applying the SOG threshold, the amount of data points significantly decrease from 7453 data points to 3506 data points. This indicated that about half of the total data points represented the ship's stationary behaviour.

From preliminary analysis, possible source of error is identified for data points representing current speed. In range of current speed between 0.01 and 0.03 [m/s], noticeable peak in data points is observed. This peak attributed to missing information on northward and eastward current speed in some data points from the provided dataset. This resulted in single random error value for current speed which resulted in the peak observed in the histogram.

To address the missing values, the missing values for eastward current and northward current are imputed using KNNImputer feature from Scikit-Learn. Each sample's missing values are imputed using the mean of nearest neighbour found in training dataset [Fabian Pedregosa et al. \(2011\)](#). Once the missing values of northward and southward current are imputed, the current speed for the missing values will be recalculated.

The imputing approach using k-nearest neighbour is also applied to other weather features that contained missing values i.e. NaN values. Imputing missing values



**Figure 3.3:** Journey of the ship in a year

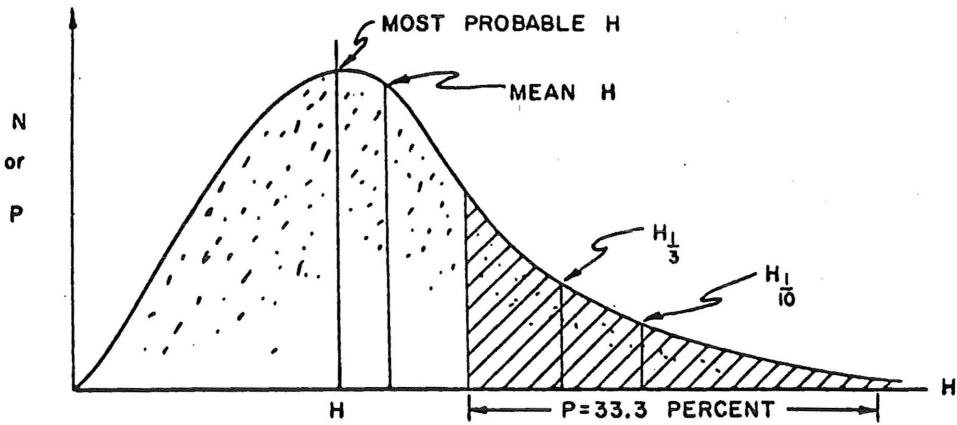
is necessary as modelling package by Scikit-Learn cannot handle missing values. Imputing strategy using k-nearest neighbour is considered as it should reflect the weather conditions within the region of missing values.

### 3.2.2 Feature Selection

To select appropriate features for the model, correlation between the features is first studied. Feature selection is necessary to simplify the model and subsequently save computing cost during training. Selection of features is based on statistical approach of High Correlation Filter proposed by Abebe et al. [Abebe et al. \(2020\)](#). This approach considers pairs of features with correlation features higher than 0.7 as one entity. However, the selection of highly correlated features must not violate natural state of matter. Therefore, in addition to statistical approach, the scientific reasoning behind the correlations will be considered and prioritised over the statistical approach.

From AIS data, the information on (1) time, (2) latitude, (3) longitude, (4) width, and (5) length are not included for training. As time, latitude and longitude have no impact on the ship. While the width and length is properties from the ship that remain constant.

The features (1) combined wind wave swell height, (2) swell height, maximum wave height (3) and wind wave height are physically correlated. In sea wave theory, wind wave swell height is also known as significant wave height  $H_{1/3}$ . It is defined as the mean of the highest one-third of waves in the wave record [Holthuijsen \(2007\)](#).



**Figure 3.4:** Statistical distribution of wave heights [Bretschneider \(1965\)](#)

The distribution of wave heights can be represented by probability density function. Hence, the term “highest one-third of waves” here means the region of wave heights that belong in the upper one-third of a probability density function, this is illustrated in Figure 3.4. From this distribution, the relation between significant wave height  $H_{1/3}$ , the highest ten percent of waves  $H_{10}$  and average wave height  $\bar{H}$  can be summarised as follows [Bretschneider \(1965\)](#); [Holthuijsen \(2007\)](#):

$$\bar{H} = 0.625 \cdot H_{1/3} \quad (3.2.1)$$

$$H_{10} = 2.03 \cdot \bar{H} = 1.27 \cdot H_{1/3} \quad (3.2.2)$$

$$H_{\max} = 2 \cdot H_{1/3} \quad (3.2.3)$$

Additionally, Bitner-Gregersen [Bitner-Gregersen \(2005\)](#) described the relation between the significant wave height, wind wave height and swell height through following equation:

$$H_{1/3} = \sqrt{(H_{\text{swell}})^2 + (H_{\text{windwave}})^2} \quad (3.2.4)$$

From here, it is clear that significant wave height should be retained for modelling, as it holds critical information regarding wave properties. The features swell height, wind wave height and maximum wave height will be dropped as it can be defined through correlations defined in Equation (3.2.1), Equation (3.2.2), Equation (3.2.3) and Equation (3.2.4). This decision is also statistically supported through the high correlation filter method. As shown in Figure 3.13, high correlation are observed between these features.

From Figure 3.13, high correlation is observed between wave period, swell period and wind wave period. Bitner-Gregersen further elaborated that the state of the sea can be described through the significant height  $H_{1/3}$  and spectral peak  $T_p$  with help of Torsethaugen peak [K. Torsethaugen et al. \(2004\)](#). Hence, the features swell

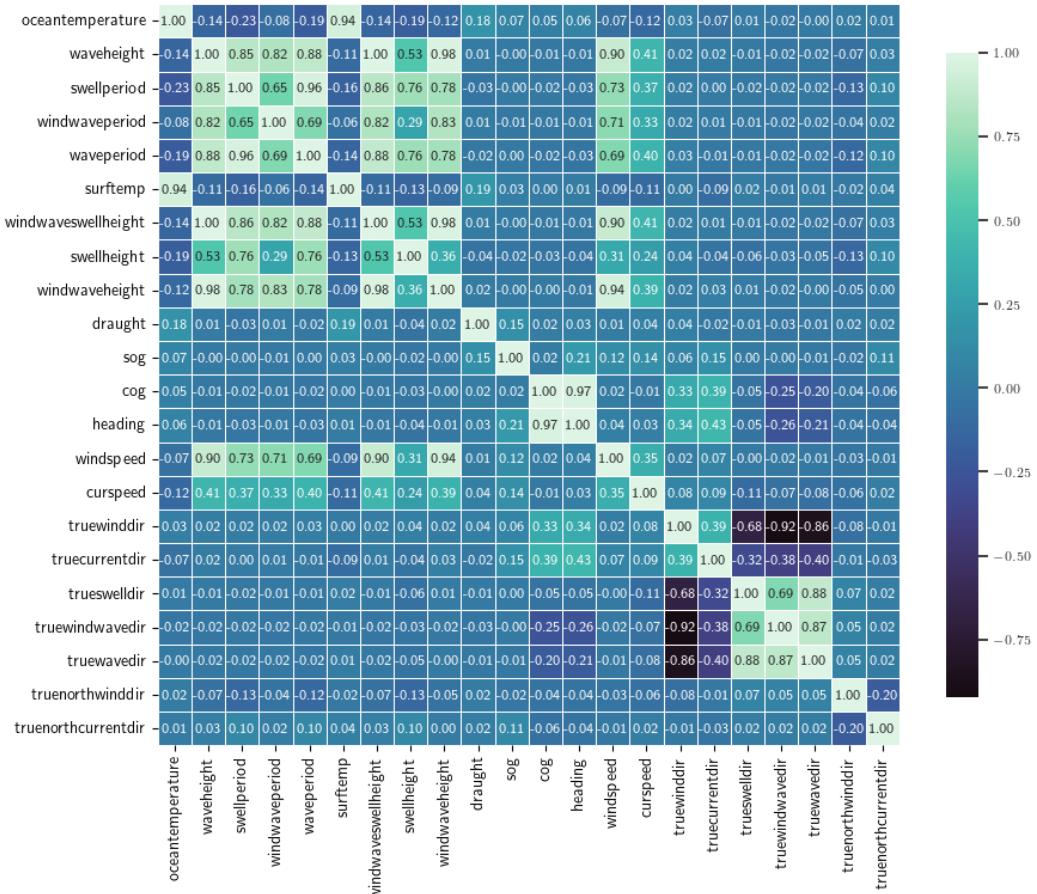


Figure 3.5: Correlation Heat Map

<b>Training Label</b>	
SOG [Knots]	sog
<b>Training Features</b>	
COG [m/s]	cog
Heading [°]	heading
Draught [m]	draught
Wind Speed [m/s]	windspeed
Air Temperature Above Oceans [K]	oceantemperature
Maximum Wave Height [m]	waveheight
Wave Period [s]	waveperiod
Sea Surface Temperature [K]	surftemp
Combined Wind Wave Swell Height [m]	windwaveswellheight
Current Speed [m/s]	curspeed
True Wind Direction [°]	truewinddir
True Current Direction [°]	truecurrentdir
True Wave Direction [°]	truelavedir

**Table 3.2:** Structure of fused dataset

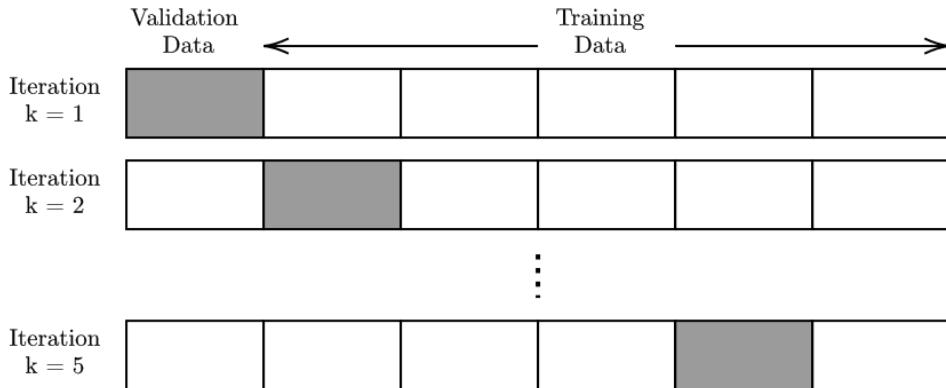
period and wind wave period are discarded as it only distinguish whether the sea is dominated by swell or by wind. The feature wave period will still be retained. Consequently, the features true wind wave direction and true swell direction will be discarded as the features that explained the magnitude of these features are discarded.

Statistically, the heading and COG are highly correlated, but both features are retained as it explain two different parameters of the ship. Course Over Ground reflects the ship course heading while heading represented the actual heading of the ship at a particular point of time. Same principle also apply between air temperature above ocean and sea surface temperature. Air temperature above oceans represents the temperature of wind while sea surface temperature represents current temperature of current.

From feature selection, 5 features from AIS data are discarded while 11 features are removed from the weather data. To predict the ship speed, The SOG will be selected as the label to train the model. The remaining attributes will be selected as training features. This is summarised in Table 3.2.

### 3.2.3 Modelling

In this section, the modelling of ship speed through SOG using selected features will be performed using tree-based regressor model. The tree-based regressor model considered are decision tree regressor, random forest regressor and extra-tree regressor. In addition, the tree-based models are compared against multiple linear regressor to as benchmark. The methodology to develop the best model is divided into several



**Figure 3.6:** Visual illustration of k-folding, Grey shaded box represents the validation data while white box represents the training data

steps.

For training, the dataset is split into training, validation and test dataset in ratio of 73:18:9. Journey data from the month of June is arbitrarily selected as test dataset. The remaining dataset will be split into training and validation dataset in 80:20 ratio. The explanation of training process and selection of the best model is broken down into several sections. In Section 3.2.5, the tuning parameter of Scikit-Learn will be studied extensively as suitable tuning could result in improved model performance.

Appropriate statistical performance measures are applied to each model; the performance measures selected will help to evaluate how well a model is able to make generalisation on validation and test dataset. The evaluation will be cross validated in form of k-folding. The details on evaluation methodology used in this thesis will be discussed in Section 3.2.4.

### 3.2.4 Performance Metrics for Model Validation

To gain sensible estimate of model performance and how precise a model is, the model will be cross validated by means of k-folding. K-fold cross validation split the training set into  $k$  subsets which is called *folds*, then the model will be trained  $k$  times using  $k-1$  subsets and remaining one for validation, this process is illustrated in Figure 3.6. For each iteration, each model is evaluated using different performance metrics such as (1) Coefficient of Determination ( $R^2$ ), (2) Explained Variance (EV), (3) Mean Absolute Error (MAE), (4) Root Mean Square (RMSE) and (5) Median Absolute Deviation (MAD). The results from each iteration is then averaged, where the information on model precision can be gained from the standard deviation. Performing k-fold cross validation checks model robustness against different datasets. The properties of each performance metric will be discussed in the following sections.

### 3.2.4.1 Coefficient of Determination ( $R^2$ )

The coefficient of determination  $R^2$  gives a measure on prediction quality,  $R^2$  quantifies the ability of the regression model to approximate the actual values.  $R^2$  is defined by Equation (3.2.5), where  $y$  represents true target output,  $\hat{y}$  represents the predictor output and  $\bar{y}$  represents the mean.  $R^2$  score range between 0 and 1, higher values i.e.  $R^2 \rightarrow 1$  indicate better model fit and score of 1 indicate perfect prediction.

$$R^2(y, \hat{y}) = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2} \quad \text{where } \bar{y} = \frac{1}{n} \sum_1^n y_i \quad (3.2.5)$$

### 3.2.4.2 Explained Variance (EV)

Explained variance indicate how well a model can capture variance from a dataset. It is defined by Equation (3.2.6), where  $\sigma_x$  represents standard deviation of parameter  $x$ . EV score range between 0 and 1, where the best score of  $EV = 1$  can be obtained if  $\sigma_{(y-\hat{y})}^2 \rightarrow 0$ .

$$EV(y, \hat{y}) = 1 - \frac{\sigma_{(y-\hat{y})}^2}{\sigma_y^2} \quad (3.2.6)$$

### 3.2.4.3 Mean Absolute Error (MAE)

MAE indicated the expected value of absolute ( $L^1$  norm) error, and it can be calculated by:

$$MAE(y, \hat{y}) = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i| \quad (3.2.7)$$

### 3.2.4.4 Root Mean Square Error (RMSE)

The RMSE describe the expected value of quadratic error. RMSE place large penalty on large deviation between true and estimated values and for this reason, it can be used to as a metric to indicate model performance against outliers. Ideal score is observed when  $RMSE \rightarrow 0$ . RMSE can be considered as absolute measure of model fitness. Omitting the root term, RMSE becomes MSE, which is the loss function of Equation (2.2.2) that is used to determine the most optimal split in a regression decision tree.

$$RMSE(y, \hat{y}) = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (3.2.8)$$

### 3.2.4.5 Median Absolute Deviation (MAD)

MAD is a performance metrics that considers the median of the absolute errors. It is robust to outlier as it only consider median performance

$$MAD(y, \hat{y}) = \text{median}(|y_1 - \hat{y}_1|, \dots, |y_i - \hat{y}_i|) \quad (3.2.9)$$

## 3.2.5 Model Hyperparameter Optimisation

The subject of parameter tuning was briefly discussed in Section 2.2.1. In Section 2.2.1 parameter tuning was applied to decision tree regressor to avoid overfitting by changing the minimum amount of samples a leaf node has. This example implies that altering model hyperparameter will affect the model performance. However, the optimisation of the hyperparameter cannot be performed *a priori* and as such iterative process will be performed until best hyperparameter value is found.

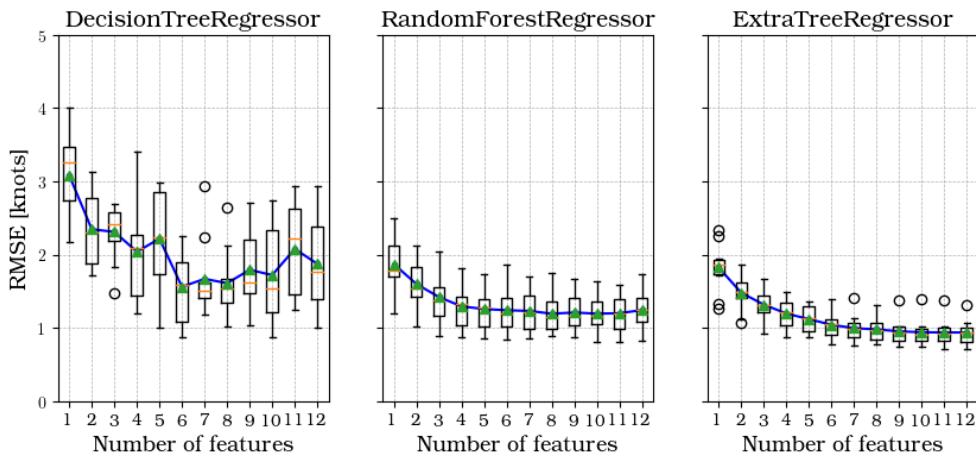
Scikit-Learn offers `GridSearchCV` and `RandomizedSearchCV` to help search for the most optimal hyperparameter. Both solutions operate with similar principle: The selected hyperparameters to be tuned with its value range is evaluated using cross validation to evaluate the best possible combination between the selected hyperparameters. The difference between `GridSearchCV` and `RandomizedSearchCV` lies in how it searches for the best value for the selected hyperparameters: `GridSearchCV` involves construction of grids containing all possible combinations of hyperparameter value in specified range. `RandomizedSearchCV` randomly samples hyperparameter values.

The exhaustive nature of `GridSearchCV` means that it is computationally costly to perform, especially when there are multiple hyperparameters to be considered and value search space is large. `RandomizedSearchCV` gives more control to computing budget by setting the number of iteration and usually produces more accurate results than `GridSearchCV` approach. [Géron \(2019\)](#); [Bergstra and Bengio \(2012\)](#).

For this reason, the `RandomizedSearchCV` will be employed to search for best possible hyperparameter. However, the limitation of *a priori* knowledge of hyperparameter value still exists. In spite of `RandomizedSearchCV` ability to control the computational budget, it is still takes considerable time to obtain the best hyperparameter value. The computational budget may be spent on searches in unpromising search space. With that, initial exploration on the effect of each hyperparameter on model performance will be performed to give better overview on which search space that should be considered during hyperparameter optimisation. In the next subsections, the effect of tunable hyperparameter of tree-based model from Scikit-Learn will be explored to give baseline numbers for the search space. RMSE is used as performance metrics as the hyperparameter parameter optimisation done in this thesis aims to reduce the error during prediction.

### 3.2.5.1 Number of features

Defined with default value as `max_features=None` in Scikit-Learn. This hyperparameter controls the number of features to be considered when looking for the best split, the default `None` option means it will consider all features. This parameter tuning is available for Decision Tree Regressor, Random Forest Regressor and ExtraTree Regressor. Initial exploration indicated Random Forest Regressor and Extra Tree Regressor benefit from considering more features, Decision Tree Regressor requires further fine-tuning to optimise the model as the default `None` means it will consider all features when searching for best split.



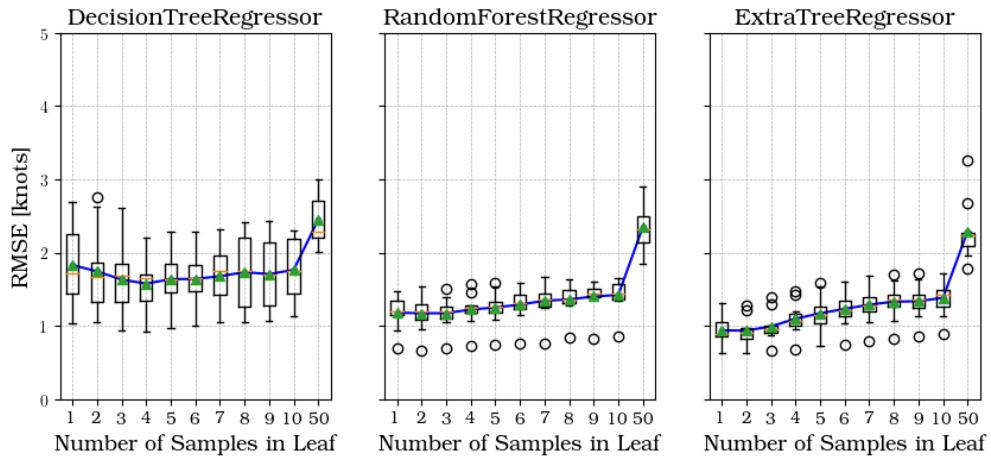
**Figure 3.7:** Hyperparameter tuning of `max_features`

### 3.2.5.2 Number of sample in a leaf node

Defined with default value as `min_samples_leaf=1` in Scikit-Learn. This parameter controls number of samples required to be at leaf node, where split point will be considered if the leaf contains at least `min_samples_leaf=n` training samples in each left and right branch. As shown in Figure 2.4, tuning this hyperparameter to higher values helps to smoothen the model and avoid overfitting. However, this may lead to underfitting as the model is unable to capture the trend within the data. This is supported by the findings shown in Figure 3.8, the DTR benefits from regularisation at certain breakeven point, in this case, it is found to be at `min_samples_leaf=4`. But after this breakeven point, the model's performance degrades. It is also observed that RFR and ETR does not benefit from any form of regularisation.

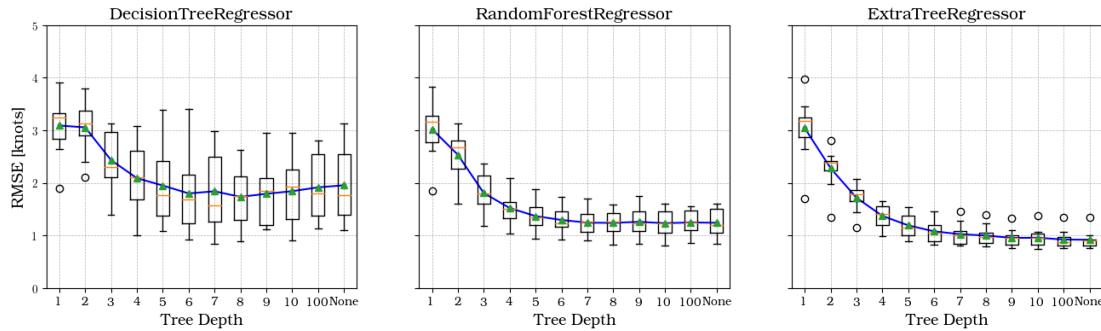
### 3.2.5.3 Depth of Tree

Defined with default value as `max_depth=None` in Scikit-Learn. This hyperparameter controls the growth of the tree. Leaving it at `max_depth=None` means the tree will grow until all leaves are pure i.e. until minimum MSE is obtained or when the number of samples is less than the minimum number of samples required to split



**Figure 3.8:** Hyperparameter tuning of `min_samples_leaf`

an internal node. Similar to `min_samples_leaf`, DTR shows improvement until a certain breakeven point. RFR performance seems to stabilise at certain depth while ETR benefits from allowing full growth of the tree. The results are summarised in Figure 3.9



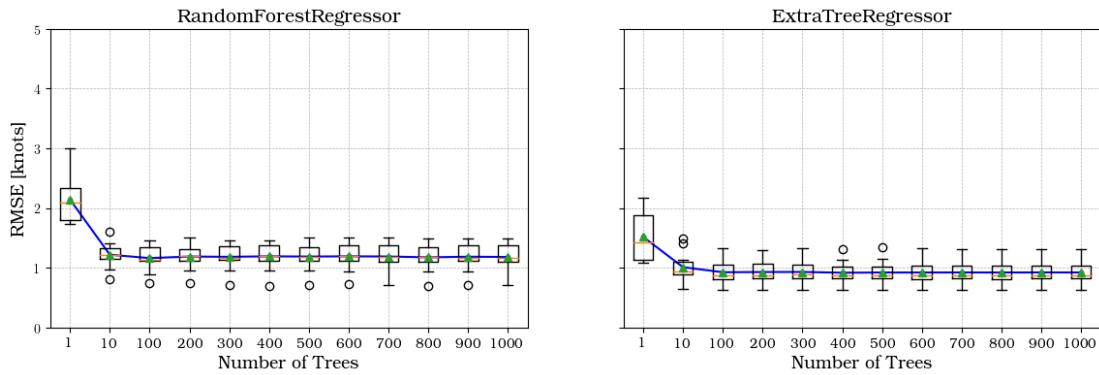
**Figure 3.9:** Hyperparameter tuning of `max_depth`

### 3.2.5.4 Number of Trees

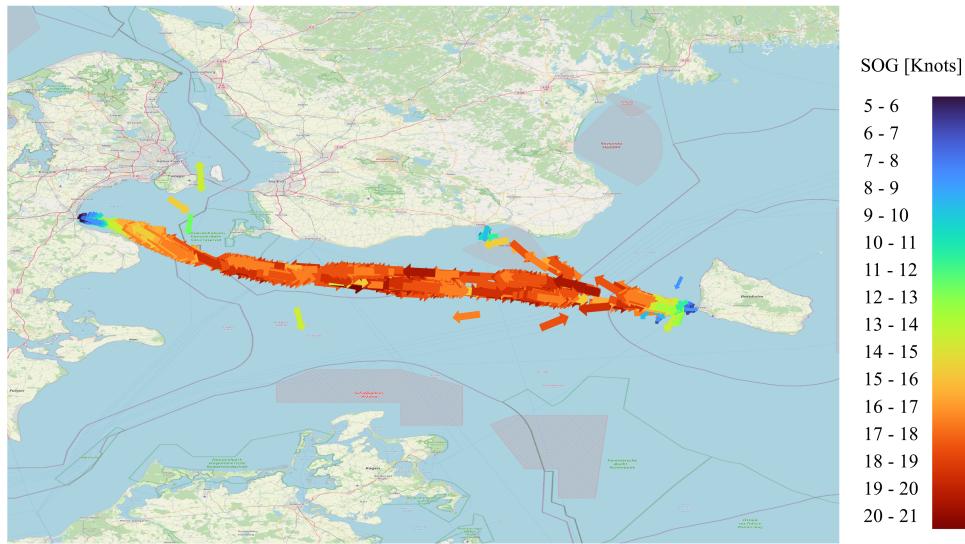
Defined with default value as `n_estimators=100`. This hyperparameter controls the amount of trees i.e. predictors in a forest. Tuning of number of trees will have an effect on the training time and it is only available to RFR and ETR. The default value seems to yield satisfactory result, as the performance for both RFR and ETR stabilise after in this case stabilise after 100 trees, as seen in Figure 3.7.

## 3.2.6 Methodology Application

- Two data sources are imported. `AIS_weather_H_ok2_copy.csv` and `AIS_weather_h_rename_copy.csv`. The information from the latter comma delimited file will be used for calculating the ship Speed Through Water (STW).



**Figure 3.10:** Hyperparameter tuning of n\_estimators



**Figure 3.11:** Journey of the ship in June

The information required is the true north current direction. Which is obtained from the vector component of the Northward and Southward current.

- This dataframe will be merged with the main dataframe from the file AIS\_weather\_H\_ok2\_copy.csv.
- Omission of the journey data between Ronne and Sassnitz
- SOG threshold is applied to omit ship mooring and maneuvering to accurately represent the ship's steady state operation [Abebe et al. \(2020\)](#); [Bal Beşikçi et al. \(2016\)](#); [Gkerekos et al. \(2019\)](#); [Yang et al. \(2020\)](#). This threshold is selected as 5 knots according to [Abebe et al. \(2020\)](#)
- The AIS data from June is filtered. This data will be used as validation data to check the model's performance.

### 3.2.7 Data Analysis

- The features are represented in a histogram plot. For the feature Current speed, anomaly is detected. Certain spike is detected around 0.01 – 0.03 m/s. Reasons unknown. The data is retained, including the spike, until a definitive answer can be found.
- OPEN QUESTION : What is the necessity of feature standardization / normalization ? Normalization is required for ANN as model training requires the value between 0 and 1. But in case of RFR, there is no such requirement. Through testing, data standardization also does not seem to improve the model's performance.

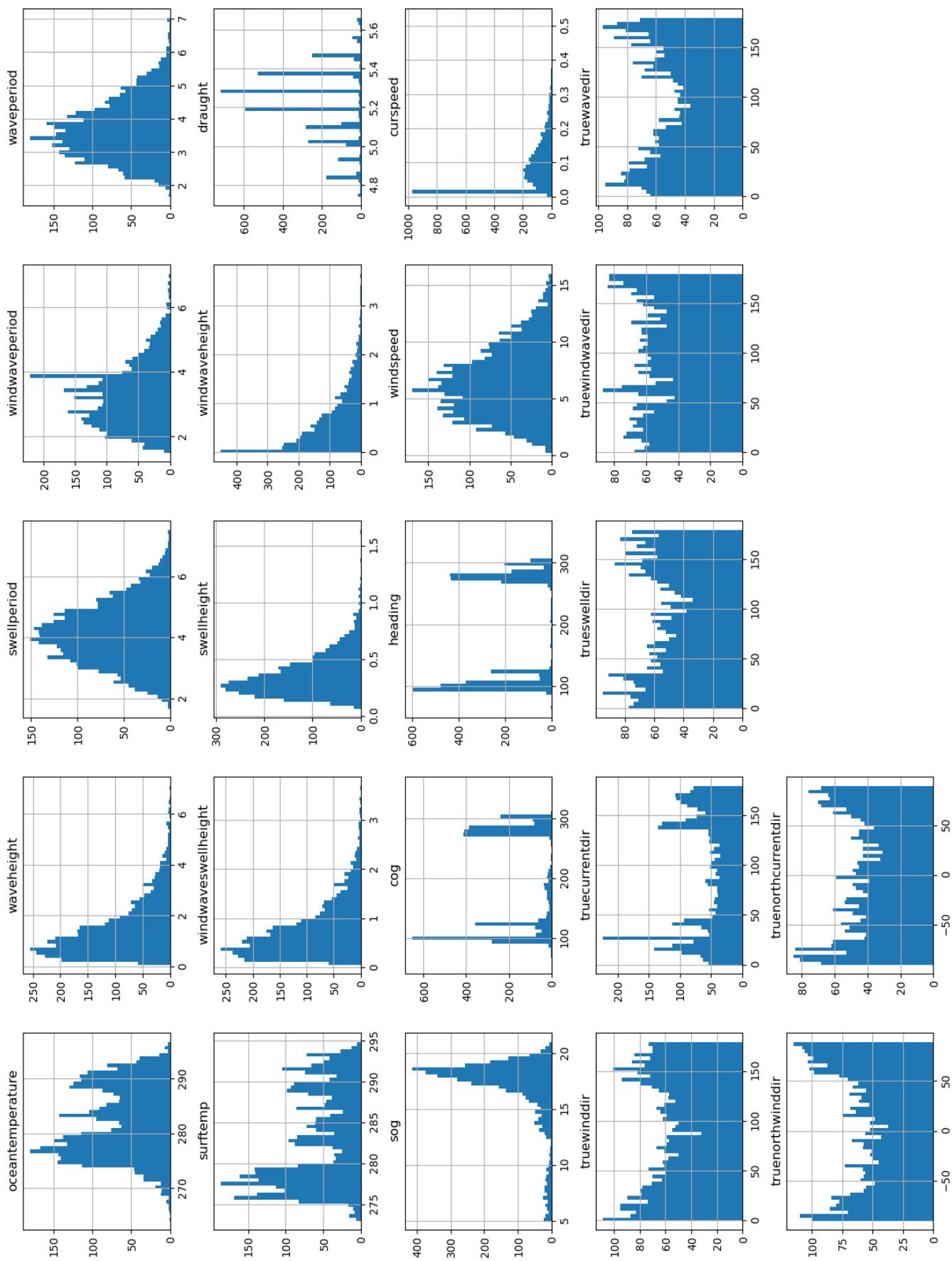


Figure 3.12: Histogram of the features

- The correlation of the features against SOG are determined. It is found that :
  - Draught
  - Course Over Ground (COG)
  - heading
  - Wind Speed
  - Current Speed
  - True Current direction

Have relatively stronger correlation to SOG compared to other features, albeit the correlation is a weak one

- The correlation between the features is displayed using the following the heat map. From the heat map it can be observed that between these features:
  - Waveheight and wind wave swell height
  - Waveheight and wind wave height
  - Windwaveswellheight and wave period

Have a strong correlation between each other.

- Open topic:
  - Feature reduction is possible, [Abebe et al. \(2020\)](#) suggested high feature correlation filter, the filter suggest that two features which has a high correlation ( $> 90\%$ ) is to be combined into a single feature. But the author is unsure whether this combination is physically sensible. Hence, this filter is yet to be applied for feature reduction.
  - Some of these features can be connected through wave equations, but the author has not found an equation which could relate these features.
- The random forest regressor could not function when NaN values are present. With that, the missing values are filled in using the `imputer` function. The missing values are filled in by means of KNN.

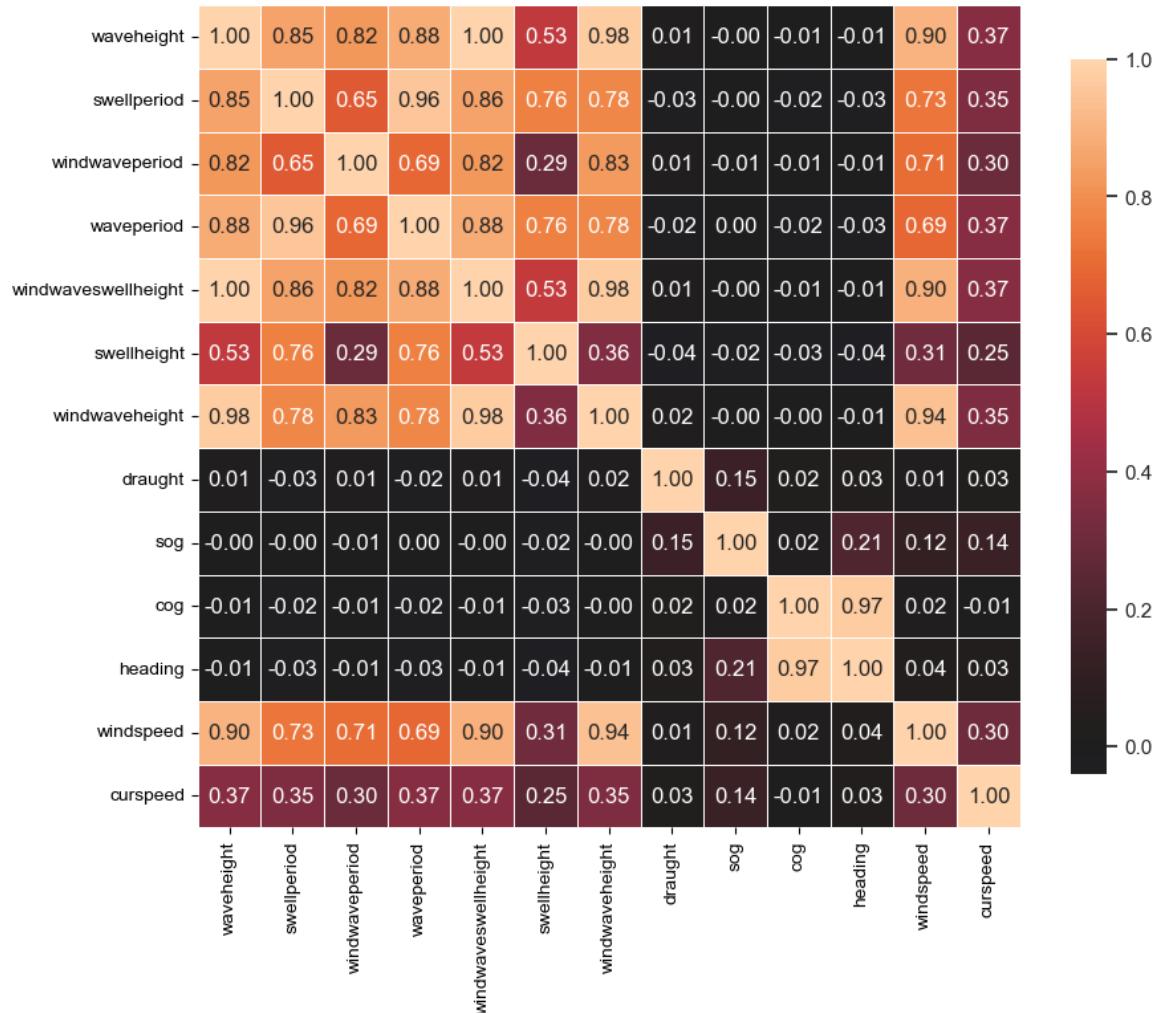


Figure 3.13: Correlation Heat Map

### 3.2.8 Modelling

- The data is split into 80:20 ratio. But considering the validation data, it is split into approximately 73:18:9.
- The model is then trained using Random Forest Regression (RFR). Additional training is also performed using Decision Tree Regressor (DTR). DTR model performance will be used as a benchmark as it is also a tree-based modelling method with similar methodology to RFR.
- The computational time of DTR is significantly faster than RFR Model Evaluation

### 3.2.9 Predicting STW

- The ship's Speed Through Water STW can be calculated using vector component of the SOG and current speed. The direction used will be according to

True North. [Yang et al. \(2020\)](#); [Zhou et al. \(2020\)](#)

- SOG represents the speed of the ship with reference to the ground, while the STW represent the ship's speed with reference to water.
- SOG also can be termed by the ship's speed that is captured by the GPS, and does not consider any effect of the current
- This means that the ship's STW will be greater than the ship's SOG when there is current moving against the ship's movement direction and vice versa
- The vector decomposition can be defined from the following equations, which is based on the equation by [Yang et al. \(2020\)](#):
  - The ship's SOG  $V_g$  can be decomposed into  $V_g^x$  and  $V_g^y$ , which represents the  $x$  and  $y$  components of the SOG respectively using the ship's course heading (COG)  $\beta$  *with respect to True North*:

$$V_g^x = V_g \sin(\beta) \quad (3.2.10)$$

$$V_g^y = V_g \cos(\beta) \quad (3.2.11)$$

- To consider the effect of sea current. The current speed  $V_c$  will also be decomposed to  $x$  and  $y$  components respectively using the current direction  $\gamma$  *with respect to True North*:

$$V_c^x = V_g \sin(\gamma) \quad (3.2.12)$$

$$V_c^y = V_g \cos(\gamma) \quad (3.2.13)$$

- from here the ship' STW  $V_{wx}$  and  $V_{wy}$  component can be found from the following equation:

$$V_w^x = V_g^x - V_c^x \quad (3.2.14)$$

$$V_w^y = V_g^y - V_c^y \quad (3.2.15)$$

- The magnitude of the STW can be readily obtained from the following vector synthesis

$$V_w = \sqrt{(V_w^x)^2 + (V_w^y)^2} \quad (3.2.16)$$

- This principle is applied to the following Python script. 3.2.12

```
1      # Convert SOG from [Knots] to [m/s]
2
3      dfprog["vgms"] = dfprog["sog_pred"]/1.9438
4
5      # Convert the angles from [Degrees] to [Radians]
6
7      rad_gamma = np.deg2rad(dfprog["gamma"])
8      rad_cog = np.deg2rad(dfprog["cog"])
9
10     # Decomposition in x-component
11
12     dfprog["vgx"] = dfprog["vgms"] * np.sin(rad_cog)
13     dfprog["vcx"] = dfprog["curspeed"] * np.sin(rad_gamma)
14     dfprog["stw_x"] = (dfprog["vgx"] - dfprog["vcx"])
15
16     # Decomposition in y-component
17
18     dfprog["vgy"] = dfprog["vgms"] * np.cos(rad_cog)
19     dfprog["vcy"] = dfprog["curspeed"] * np.cos(rad_gamma)
20     dfprog["stw_y"] = (dfprog["vgy"] - dfprog["vcy"])
21
22     # Vector synthesis and reconversion to [Knots] from [m/s]
23
24     dfprog["vwms_p"] = np.sqrt(dfprog["stw_x"]**2 + dfprog["stw_y"]**2)
25     dfprog["stw_pred"] = dfprog["vwms_p"]*1.9438
26
27
28
```

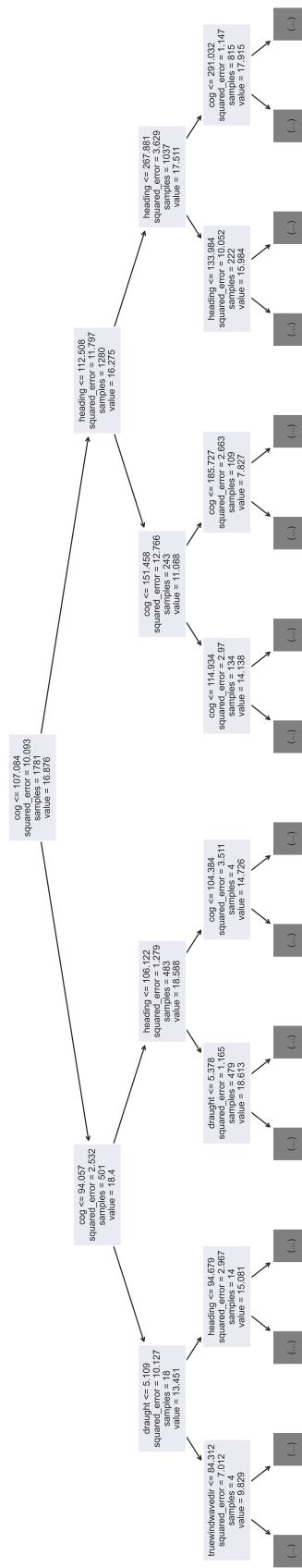


Figure 3.14: Correlation Heat Map

# Chapter 4

## Result and Discussion

The result of the research is discussed in this chapter. This comprises model validation and how different statistical metrics are used to analyze the model's performance.

### 4.0.1 Model Evaluation

The model are tested against four metrics, namely:

- $R^2$  : Indicate model fit. Best Score = 1
- Explained Variance EV : Indicate amount of variance in model. Best Score = 1
- Mean Absolute Error MAE : Indicate how much error a model makes in its prediction. Best Score = 0
- Root Mean Square Error RMSE : Same as MAE, more sensitive to outlier. Best Score = 0
- Median Absolute Error MAD : Check robustness against outlier. Best Score = 1

The result is summarized in the following table

Model	RFR	DTR	LR
$R^2$	0.9328181446941499	0.8526085810220092	1
EV	0.932872958708872	0.8526260247615258	2
MAE	0.5546347329650284	0.8108982427834758	3
RMSE	0.7095480848510665	1.5566896535262504	4
MAD	0.38484635910000087	0.5475717149999983	5

Table 4.1: Model performance

Model	RFR	DTR	LR
$R^2$	0.9328181446941499	0.8526085810220092	1
EV	0.932872958708872	0.8526260247615258	2
MAE	0.5546347329650284	0.8108982427834758	3
RMSE	0.7095480848510665	1.5566896535262504	4
MAD	0.38484635910000087	0.5475717149999983	5

**Table 4.2:** Model performance

# **Chapter 5**

## **Summary and Outlook**

In this chapter the summary of this research will be discussed. This section includes reflections of the research process and presents any possible suggestions and recommendations in this line of research. This chapter concludes this thesis.

# Bibliography

- Misganaw Abebe, Yongwoo Shin, Yoojeong Noh, Sangbong Lee, and Inwon Lee. Machine learning approaches for ship speed prediction towards energy efficient shipping. *Applied Sciences*, 10(7):2325, 2020. doi:[10.3390/app10072325](https://doi.org/10.3390/app10072325).
- E. Bal Beşikçi, O. Arslan, O. Turan, and A. I. Ölcer. An artificial neural network based decision support system for energy efficient ship operations. *Computers & Operations Research*, 66:393–401, 2016. ISSN 03050548. doi:[10.1016/j.cor.2015.04.004](https://doi.org/10.1016/j.cor.2015.04.004).
- J Bergstra and Yoshua Bengio. Random search for hyper-parameter optimization. *J. Mach. Learn. Res.*, 2012. URL <https://www.semanticscholar.org/paper/Random-Search-for-Hyper-Parameter-Optimization-Bergstra-Bengio/188e247506ad992b8bc62d6c74789e89891a984f>.
- Nicolas Bialystocki and Dimitris Konovessis. On the estimation of ship's fuel consumption and speed curve: A statistical approach. *Journal of Ocean Engineering and Science*, 1(2):157–166, 2016. ISSN 24680133. doi:[10.1016/j.joes.2016.02.001](https://doi.org/10.1016/j.joes.2016.02.001).
- Elzbieta M. Bitner-Gregersen. Joint probabilistic description for combined seas. In *24th International Conference on Offshore Mechanics and Arctic Engineering: Volume 2*, pages 169–180. ASMEDC, 2005. ISBN 0-7918-4196-0. doi:[10.1115/OMAE2005-67382](https://doi.org/10.1115/OMAE2005-67382).
- Leo Breiman. Bagging predictors. *Machine Learning*, 24(2):123–140, 1996. ISSN 1573-0565. doi:[10.1007/BF00058655](https://doi.org/10.1007/BF00058655). URL <https://link.springer.com/article/10.1007/bf00058655>.
- Leo Breiman. Random forests. *Machine Learning*, 45(1):5–32, 2001. ISSN 1573-0565. doi:[10.1023/A:1010933404324](https://doi.org/10.1023/A:1010933404324). URL <https://link.springer.com/article/10.1023/a:1010933404324>.
- Leo Breiman, Jerome H. Friedman, Richard A. Olshen, and Charles J. Stone. *Classification And Regression Trees*. Routledge, 2017. ISBN 9781315139470. doi:[10.1201/9781315139470](https://doi.org/10.1201/9781315139470). URL <https://www.taylorfrancis.com/books/mono/10.1201/9781315139470/classification-regression-trees-leo-breiman>.
- Charles L. Bretschneider. *Generation of waves by wind. State of the art.* 1965. URL <https://apps.dtic.mil/sti/citations/ad0612006>.

Andrea Coraddu, Luca Oneto, Francesco Baldi, and Davide Anguita. Vessels fuel consumption forecast and trim optimisation: A data analytics perspective. *Ocean Engineering*, 130:351–370, 2017. ISSN 00298018.

doi:[10.1016/j.oceaneng.2016.11.058](https://doi.org/10.1016/j.oceaneng.2016.11.058). URL

<https://www.sciencedirect.com/science/article/pii/S0029801816305571>.

Danish Maritime Authority. Safety at sea, navigational information, AIS data, 2023.

URL <https://dma.dk/safety-at-sea/navigational-information/ais-data>.

Yuquan Du, Qiang Meng, Shuaian Wang, and Haibo Kuang. Two-phase optimal solutions for ship speed and trim optimization over a voyage using voyage report data. *Transportation Research Part B: Methodological*, 122:88–114, 2019. ISSN 01912615. doi:[10.1016/j.trb.2019.02.004](https://doi.org/10.1016/j.trb.2019.02.004).

Fabian Pedregosa, Gaël Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas, Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, and Édouard Duchesnay. Scikit-learn: Machine learning in python. *Journal of Machine Learning Research*, 12(85):2825–2830, 2011. URL <http://jmlr.org/papers/v12/pedregosa11a.html>.

Aurélien Géron. *Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow: Concepts, tools, and techniques to build intelligent systems* / Aurélien Géron. O'Reilly, Sebastopol, CA, second edition edition, 2019. ISBN 978-1-492-03264-9.

Pierre Geurts, Damien Ernst, and Louis Wehenkel. Extremely randomized trees. *Machine Learning*, 63(1):3–42, 2006. ISSN 1573-0565.

doi:[10.1007/s10994-006-6226-1](https://doi.org/10.1007/s10994-006-6226-1). URL

<https://link.springer.com/article/10.1007/s10994-006-6226-1>.

Christos Gkerekos, Iraklis Lazakis, and Gerasimos Theotokatos. Machine learning models for predicting ship main engine fuel oil consumption: A comparative study. *Ocean Engineering*, 188:106282, 2019. ISSN 00298018.

doi:[10.1016/j.oceaneng.2019.106282](https://doi.org/10.1016/j.oceaneng.2019.106282).

Alon Halevy, Peter Norvig, and Fernando Pereira. The unreasonable effectiveness of data. *IEEE Intelligent Systems*, 24(2):8–12, 2009. ISSN 1941-1294.

doi:[10.1109/MIS.2009.36](https://doi.org/10.1109/MIS.2009.36).

Michael Haranen, Pekka Pakkanen, Risto Kariranta, and Jouni Salo. White, grey and black-box modelling in ship performance evaluation. 2016. URL [https://www.researchgate.net/publication/301355727\\_White\\_Grey\\_and\\_Black-Box\\_Modelling\\_in\\_Ship\\_Performance\\_Evaluation](https://www.researchgate.net/publication/301355727_White_Grey_and_Black-Box_Modelling_in_Ship_Performance_Evaluation).

Trevor Hastie, Robert Tibshirani, and J. H. Friedman. *The elements of statistical learning: Data mining, inference, and prediction* / Trevor Hastie, Robert Tibshirani, Jerome Friedman. Springer series in statistics. Springer, New York, 2nd ed. edition, 2009. ISBN 9780387848570. doi:[10.1007/b94608](https://doi.org/10.1007/b94608).

- Leo H. Holthuijsen. *Waves in oceanic and coastal waters*. Cambridge University Press, Cambridge, 2007. ISBN 9780521860284.
- J. Holtrop. A statistical re-analysis of resistance and propulsion data. *Published in International Shipbuilding Progress, ISP, Volume 31, Number 363, 1984*. URL <https://repository.tudelft.nl/islandora/object/uuid%3Aca12a502-fc85-45e4-a078-db7284127e3c>.
- J. Holtrop and G.G.J. Mennen. A statistical power prediction method. *Netherlands Ship Model Basin, NSMB, Wageningen, Publication No. 603, Published in: International Shipbuilding Progress, ISP, Volume 25, Number 290, October 1978, 1978*. URL <https://repository.tudelft.nl/islandora/object/uuid%3A62c40df8-18cc-4225-a65a-54ff5c1609fb>.
- J. Holtrop and G.G.J. Mennen. An approximate power prediction method. *Netherlands Ship Model Basin, NSMB, Wageningen, Publication No. 689, Published in: International Shipbuilding Progress, ISP, Volume 29, Nr 335, 1982, 1982*. URL <https://repository.tudelft.nl/islandora/object/uuid%3Aee370fed-4b4f-4a70-af77-e14c3e692fd4>.
- IMO. Revised guidelines for the onboard operational use of shipborne Automatic Identification Systems (AIS), 2015. URL <https://www.imo.org/en/OurWork/Safety/Pages/AIS.aspx>.
- IMO. Fourth IMO GHG study 2020. *International Maritime Organization London, UK, 2020*.
- Jóan Petur Petersen. Mining of ship operation data for energy conservation. 2011. URL <https://orbit.dtu.dk/en/publications/mining-of-ship-operation-data-for-energy-conservation>.
- K. Torsethaugen, S. Haver, and S. Norway. Simplified double peak spectral model for ocean waves. 2004. URL <https://www.semanticscholar.org/paper/Simplified-Double-Peak-Spectral-Model-For-Ocean-Torsethaugen-Haver/0f1b1509791d441861ff6c2940dd13b1f939f149>.
- Seong-Hoon Kim, Myung-Il Roh, Min-Jae Oh, Sung-Woo Park, and In-Il Kim. Estimation of ship operational efficiency from ais data using big data technology. *International Journal of Naval Architecture and Ocean Engineering*, 12:440–454, 2020. ISSN 2092-6782. doi:[10.1016/j.ijnaoe.2020.03.007](https://doi.org/10.1016/j.ijnaoe.2020.03.007). URL <https://www.sciencedirect.com/science/article/pii/S2092678220300091>.
- Max Kuhn and Kjell Johnson. *Applied predictive modeling*. Springer, New York, 2013. ISBN 9781461468486.
- Xiaohe Li, Yuquan Du, Yanyu Chen, Son Nguyen, Wei Zhang, Alessandro Schönborn, and Zhuo Sun. Data fusion and machine learning for ship fuel efficiency modeling: Part i – voyage report data and meteorological data. *Communications in Transportation Research*, 2:100074, 2022. ISSN 27724247. doi:[10.1016/j.commtr.2022.100074](https://doi.org/10.1016/j.commtr.2022.100074).

- Joan P. Petersen, Ole Winther, and Daniel J. Jacobsen. A machine-learning approach to predict main energy consumption under realistic operational conditions. *Ship Technology Research*, 59(1):64–72, 2012a. ISSN 0937-7255. doi:[10.1179/str.2012.59.1.007](https://doi.org/10.1179/str.2012.59.1.007).
- Jón Petur Petersen, Daniel J. Jacobsen, and Ole Winther. Statistical modelling for ship propulsion efficiency. *Journal of Marine Science and Technology*, 17(1): 30–39, 2012b. ISSN 0948-4280. doi:[10.1007/s00773-011-0151-0](https://doi.org/10.1007/s00773-011-0151-0).
- Stian Glomvik Rakke. *Ship emissions calculation from AIS*. PhD thesis, NTNU, 2016. URL <https://ntnuopen.ntnu.no/ntnu-xmlui/handle/11250/2410741>.
- D. Ronen. The effect of oil price on containership speed and fleet size. *Journal of the Operational Research Society*, 62(1):211–216, 2011. ISSN 0160-5682. doi:[10.1057/jors.2009.169](https://doi.org/10.1057/jors.2009.169).
- David Ronen. The effect of oil price on the optimal speed of ships. *The Journal of the Operational Research Society*, 33(11):1035, 1982. ISSN 01605682. doi:[10.2307/2581518](https://doi.org/10.2307/2581518).
- Smith, J. P. Jalkanen, B. A. Anderson, J. J. Corbett, J. Faber, S. Hanayama, E. O’Keefe, S. Parker, L. Johansson, L. Aldous, C. Raucci, M. Traut, S. Ettinger, D. Nelissen, D. S. Lee, S. Ng, A. Agrawal, J. J. Winebrake, M. Hoen, S. Chesworth, and A. Pandey. Third imo greenhouse gas study 2014. 2015. URL <https://research.manchester.ac.uk/en/publications/third-imo-greenhouse-gas-study-2014>.
- Omer Soner, Emre Akyuz, and Metin Celik. Use of tree based methods in ship performance monitoring under operating conditions. *Ocean Engineering*, 166: 302–310, 2018. ISSN 00298018. doi:[10.1016/j.oceaneng.2018.07.061](https://doi.org/10.1016/j.oceaneng.2018.07.061). URL <https://www.sciencedirect.com/science/article/pii/S0029801818314446>.
- Stopford. The organization of the shipping market. page 47, 2009.
- Tin Kam Ho. Random decision forests. In *Proceedings of 3rd International Conference on Document Analysis and Recognition*, pages 278–282 vol.1, 1995. doi:[10.1109/ICDAR.1995.598994](https://doi.org/10.1109/ICDAR.1995.598994).
- Shuaian Wang and Qiang Meng. Sailing speed optimization for container ships in a liner shipping network. *Transportation Research Part E: Logistics and Transportation Review*, 48(3):701–714, 2012. ISSN 13665545. doi:[10.1016/j.tre.2011.12.003](https://doi.org/10.1016/j.tre.2011.12.003). URL <https://www.sciencedirect.com/science/article/pii/S1366554511001554>.
- Yuanqiao Wen, Xiaoqiao Geng, Lichuan Wu, Tsz Leung Yip, Liang Huang, and Dingyong Wu. Green routing design in short seas. *International Journal of Shipping and Transport Logistics*, 9(3):371, 2017. ISSN 1756-6517. doi:[10.1504/IJSTL.2017.083474](https://doi.org/10.1504/IJSTL.2017.083474).

- N. Wijnolst, Tor Wergeland, and Kai Levander. *Shipping Innovation*. IOS Press, 2009. ISBN 9781586039431.
- Ran Yan, Shuaian Wang, and Yuquan Du. Development of a two-stage ship fuel consumption prediction and reduction model for a dry bulk ship. *Transportation Research Part E: Logistics and Transportation Review*, 138:101930, 2020. ISSN 13665545. doi:[10.1016/j.tre.2020.101930](https://doi.org/10.1016/j.tre.2020.101930).
- Ran Yan, Shuaian Wang, and Harilaos N. Psaraftis. Data analytics for fuel consumption management in maritime transportation: Status and perspectives. *Transportation Research Part E: Logistics and Transportation Review*, 155:102489, 2021. ISSN 13665545. doi:[10.1016/j.tre.2021.102489](https://doi.org/10.1016/j.tre.2021.102489). URL <https://www.sciencedirect.com/science/article/pii/S1366554521002519>.
- Liqian Yang, Gang Chen, Jinlou Zhao, and Niels Gorm Malý Rytter. Ship speed optimization considering ocean currents to enhance environmental sustainability in maritime shipping. *Sustainability*, 12(9):3649, 2020. doi:[10.3390/su12093649](https://doi.org/10.3390/su12093649).
- Yang Zhou, Winnie Daamen, Tiedo Vellinga, and Serge Hoogendoorn. *AIS data analysis for the impacts of wind and current on ship behavior in straight waterways*. 2017.
- Yang Zhou, Winnie Daamen, Tiedo Vellinga, and Serge P. Hoogendoorn. Impacts of wind and current on ship behavior in ports and waterways: A quantitative analysis based on ais data. *Ocean Engineering*, 213:107774, 2020. ISSN 00298018. doi:[10.1016/j.oceaneng.2020.107774](https://doi.org/10.1016/j.oceaneng.2020.107774).

## **Declaration in lieu of oath**

I hereby solemnly declare that I have independently completed this work or, in the case of group work, the part of the work that I have marked accordingly. I have not made use of the unauthorised assistance of third parties. Furthermore, I have used only the stated sources or aids and I have referenced all statements (particularly quotations) that I have adopted from the sources I have used verbatim or in essence.

I declare that the version of the work I have submitted in digital form is identical to the printed copies submitted.

I am aware that, in the case of an examination offence, the relevant assessment will be marked as ‘insufficient’ (5.0). In addition, an examination offence may be punishable as an administrative offence (Ordnungswidrigkeit) with a fine of up to €50,000. In cases of multiple or otherwise serious examination offences, I may also be removed from the register of students.

I am aware that the examiner and/or the Examination Board may use relevant software or other electronic aids in order to establish an examination offence has occurred

I solemnly declare that I have made the previous statements to the best of my knowledge and belief and that these statements are true and I have not concealed anything.

I am aware of the potential punishments for a false declaration in lieu of oath and in particular of the penalties set out in Sections 156 and 161 of the German Criminal Code (Strafgesetzbuch; StGB), which I have been specifically referred to.

### **Section 156 False declaration in lieu of an oath**

Whoever falsely makes a declaration in lieu of an oath before an authority which is competent to administer such declarations or falsely testifies whilst referring to such a declaration incurs a penalty of imprisonment for a term not exceeding three years or a fine.

### **Section 161 Negligent false oath; negligent false declaration in lieu of oath**

(1) Whoever commits one of the offences referred to in Sections 154 to 156 by negligence incurs a penalty of imprisonment for a term not exceeding one year or a fine. (2) No penalty is incurred if the offender corrects the false statement in time.

The provisions of Section 158 (2) and (3) apply accordingly.

---

Place,date

---

Signature