# Development of a two-stage ship fuel consumption prediction and reduction model for a dry bulk ship

## Abstract

Shipping industry is the backbone of global trade. However, the large quantities of greenhouse gas emissions from shipping, such as carbon dioxide ($CO_2$), cannot be ignored. In order to comply with the international environmental regulations as well as to increase commercial profits, shipping companies have stronger motivations to improve ship energy efficiency. In this study, a two-stage ship fuel consumption prediction and reduction model is proposed for a dry bulk ship. At the first stage, a fuel consumption prediction model based on random forest regressor is proposed and validated. The prediction model takes into account ship sailing speed, total cargo weight, and sea and weather conditions and then predicts the hourly fuel consumption of the main engine. The mean absolute percentage error of the random forest regressor is 7.91%. At the second stage, a speed optimization model is developed based on the prediction model proposed at the first stage while guaranteeing the estimated arrival time to the destination port. Numerical experiment on two consecutive-8-day voyages shows that the proposed model can reduce ship fuel consumption by 2% to 7%. The reduction in ship fuel consumption will also lead to lower $CO_2$ emissions.

*Keywords:* Fuel consumption prediction, Ship fuel efficiency, Ship speed optimization, Random forest regressor; Machine learning

## 1. Introduction

In the past few years, improving ship energy efficiency has received wide attention not only from governmental and non-governmental organizations, but also from shipping companies (Yang et al., 2019). Although shipping is a vital component of global economy, air pollutants and greenhouse gas emissions from shipping industry caused by fuel consumption cannot be ignored. Regarding greenhouse gas emissions, such as $CO_2$, it is reported by the International Maritime Organization (IMO) that $CO_2$ emissions from shipping constitute 3.1% of global emissions, while international shipping emissions take up for 2.6% of the global emissions during 2007 to 2012 (IMO 2014). Thus, an increasing number of international regulations have been focused on improving ship energy efficiency. The first related international regulation is the amendments of the International Convention for the Prevention of Pollution from Ships (MARPOL) Annex VI proposed by Marine Environment Protection Committee (MEPC) in 2011 (IMO 2011). In addition, an approach named "Energy Efficiency Operational Indicator (EEOI)" was also proposed as a monitoring tool to manage ship and fleet efficiency performance. In 2016, amendments to MARPOL Annex VI mandatorily required ships to record and report their fuel oil consumption.

For shipping companies, due to high fuel prices, fuel costs have become the dominant factor of ship operational costs (Du et al., 2019). It is estimated that ship fuel costs constitute 20% to 50% of the total ship operating costs (Leifsson, 2008; Hasselaar, 2011). For a large container ship, fuel costs can reach about three-quarters of its operating costs when the fuel prices are high. In addition, the costs can be higher if the container ship chooses to use cleaner fuel (Ronen, 2011). In shipping industry, slow steaming is a commonly used countermeasure to reduce fuel consumption, but on-time delivery may not be guaranteed (Lee et al., 2015). Thus, in order to conform to the international environmental protection regulations as well as to increase revenue and enhance competitiveness, shipping companies are developing stronger motivations to propose practicable measures to increase ship energy efficiency.

For the existing ships, it can be hard to change their structure to reduce fuel consumption. Thus, finely planning ship voyages, e.g., adopting weather routing and optimizing sailing speed are more popular measures. For the fixed sailing routes over a voyage, one main duty for the shipping company is to plan the daily sailing speeds of the ships in advance to minimize fuel consumption over the voyage while guaranteeing on-time arrival. Sailing speed optimization requires predicting ship fuel consumption in different situations. However, there are several challenges in making accurate prediction. First, inaccuracy exists in ship sailing data that can be used to construct fuel consumption prediction models, as these datasets mainly come from manually filled ship log data, such as noon reports. Second, factors influencing ship fuel consumption

are high-dimensional. Although it is widely believed that ship sailing speed is the most important influencing factor on ship fuel consumption (Fagerholt et al., 2010; Corbett et al, 2009; Psaraftis and Knotovas, 2013, Bialystocki and Konovessis, 2016), other factors can also have impacts. These factors include but are not limited to trim condition, displacement and draft conditions, weather and sea conditions, and hull and propeller roughness (Andersen et al., 2005; IMO, 2011; Bialystocki and Konovessis, 2016). Nevertheless, it is hard to have detailed information on all the influencing factors on fuel consumption, which prevents classic regression models from making accurate fuel consumption prediction. Third, as different ships have different properties and structures, one fuel consumption prediction model cannot be universally applied (Banawan et al., 2013). Alternatively, a tailored prediction model should be developed for each single ship to achieve more satisfactory prediction performance. Developing tailored machine learning models is a desirable and promising way to deal with these challenges. Machine learning models have the ability to handle multi-dimensional input data and to extract hidden information from complex datasets. In addition, they usually have better ability to deal with noisy data. Compared with traditional statistical regression models, machine learning models can address higher dimensional data (e.g., ship displacement conditions, sea and weather conditions, trim conditions, and sailing speed) and make much more accurate predictions, and thus provide a more reliable foundation on developing tailored ship fuel consumption reduction models.

The purpose of this study is to propose a two-stage ship sailing speed optimization model for a dry bulk ship which contains two steps: in step 1, a machine learning model performing regression task (i.e. a random forest regression model) with high accuracy is proposed to make predictions on ship fuel consumption under different sailing speeds as well as cargo, weather, and sea conditions; in step 2, a sailing speed optimization model is proposed based on the prediction results in step 1 to minimize ship total fuel consumption over a voyage.

## 2. Literature review
### 2.1 Research on ship fuel consumption prediction

During the last few years, there has been an increasing amount of literature on prediction of ship fuel consumption (Zhao and Yang, 2018; Yang et al., 2019). The pioneering, basic, and commonly used models are deterministic models, which are also called white box models. In a deterministic model, the ship behavior of hull resistance, propeller propulsion, and main engine performance are described (Yang et al., 2019). Typical and pioneering studies include Holtrop (1977, 1978), Holftrop and Mennen (1978), and modern studies include Kristensen and Lützen (2012). Apart from the deterministic model, two types of models are also widely used in more recent research:

statistical models and machine learning models. Regarding the development of statistical models for fuel consumption prediction, Bocchetti et al. (2013) proposed a multiple linear regression analysis model, which took ship sailed distance and displacement as well as wind speed conditions into account to predict fuel consumption and $CO_2$ emissions of a cruise ship. Bochetti et al. (2015) then developed another multiple linear regression model for a cruise ship by containing more influencing factors. Erto et al. (2015) also developed a multiple linear regression model for a cruise ship by taking ship operational factors and wind condition into consideration. As the foundation of a ship fuel consumption analysis system, Kee et al. (2018) proposed a multiple linear regression method to estimate fuel consumption of two tugboats. Although statistical models are intuitive and interpretable, there can be some drawbacks. First, parametric statistical models require making assumptions on data distributions before developing models, and this may bring bias. In addition, even if the log-log model can express the power function of speed and fuel consumption, the linear regression models usually cannot perform well when dealing with complicated data and multicollinearity data. Moreover, they are easily influenced by noisy data (Neter et al., 1996; Goldstein, 2011).

Over the past years, a growing body of innovative literature has focused on developing machine learning methods for ship fuel consumption prediction. The most popular method is Artificial Neural Networks (ANNs) model. Pedersen and Larsen (2009) proposed an ANN model for predicting propulsion power of a tanker based on ship noon report data. They also found that by combining sea and wind information, the performance of ANN model could be significantly improved. Beşikçi et al. (2016) developed a decision support system (DSS) for improving energy efficiency of an oil tanker. The decision system contained two parts: an ANN model for fuel consumption prediction under various operational conditions and a DSS based on the prediction results for energy-efficient ship operations. In comparison studies, they reported that the performance of the ANN model was superior to multiple regression analysis based on their dataset. Petersen and Jacobsen (2012a) compared the performance of ANN and Gaussian processes (GP) models when applied to predict fuel consumption of a domestic ferry. The result indicated that the performance of ANN was a little superior than the GP in all the tests. Petersen et al. (2012b) proposed tapped-delay neural network model for fuel consumption prediction of a tanker, which was then applied to trim optimization of the tanker. Petursson (2009) developed five machine learning models for fuel consumption prediction of a passenger ship: support vector regression (SVR), k-nearest neighbor (kNN), ANN, classification and regression trees (CART) and bagging. They found that the SVR and kNN outperformed the other models on their dataset. Other types of machine learning models are also adopted for ship fuel

consumption prediction. A least absolute shrinkage and selection operator (LASSO) regression model, which contained sea and weather conditions, was adopted to predict fuel consumption of a container ship (Wang et al., 2018). Soner et al. (2018) developed three tree-based models: bagging, random forest, and bootstrap based on the log dataset of a ferry ship that was also used by Petersen et al. (2012b). They identified that the performance of tree-based prediction models had higher prediction accuracy. Grey-box models, which is in between the white-box model and black-box model, were also developed. More specifically, one type of the grey-box model structure is built based on basic principles of ship propulsion and the unknown parameters are estimated by statistical regression models, such as Journée et al. (1987), Lu et al. (2013), Meng et al. (2016) and Yang et al. (2019). The other type of grey-box model combines white-box model, which describes some components of resistance or fuel consumption, and black-box model, such as machine learning and statistical models, for the remaining parts. This type of grey-box model can be seen in Leifsson et al. (2008), Coraddu et al. (2015), Haranen et al. (2016), and Coraddu et al. (2017). The advantage of grey-box models is they are able to integrate mechanistic knowledge with data analysis methods.

Machine learning models are capable of dealing with high-dimensional data and making more accurate predictions on complicated data than traditional regression models. In addition, no human interventions are needed when learning the models (Bishop, 2006; Alpaydin, 2009). Several studies have shown that the machine learning models outperform statistical models (Petersen and Jacobsen, 2012a; Wang et al., 2018; Du et al., 2019).

Regarding the factors that influence ship fuel consumption, almost all the above-mentioned studies, either based on statistical regression methods or machine learning methods, show that ship sailing speed is the dominant factor for ship fuel consumption prediction (Bocchetti et al., 2013, 2015; Petersen and Jacobsen, 2012a; Meng et al., 2016). Actually, the "cubic law" between ship sailing speed and fuel consumption, i.e., the bunker consumption of a ship in one time unit is proportional to the sailing speed to the power of three, is widely-believed and adopted in shipping industry and maritime studies (Meng et al., 2016). Apart from sailing speed, ship displacement, such as total weight of the ship, cargo conditions, and ballast water, can also have an influence on fuel consumption based on vessel dynamics. Sea conditions, such as ocean currents (Lo and McCord, 1995), sea waves and swell (Lu et al., 2015; MAN Diesel & Turbo, 2011), are also proved to be influential to ship fuel consumption. Moreover, weather conditions are regarded as relevant to ship fuel consumption. For example, Kwon (1981) and Townsin and Kwon (1993) investigated weather conditions on ship performance and a group of regression models were proposed. Recently, models incorporating sea and weather information, including wind direction and force, sea wave direction and height,

and sea water temperature have exhibited high accuracy in ship fuel consumption prediction, such as the models proposed by Wang et al. (2016), Lee et al. (2018), Meng et al. (2016), and Du et al. (2019). Combining ship sailing related features together with sea and weather conditions have shown great potential in ship fuel consumption prediction and management.

## 2.2 Research on improving ship energy efficiency

Much of the current literature on ship energy efficiency pays particular attention to finding viable measures to reduce ship fuel consumption. As suggested by SEEMP, there are several effective ways to save ship fuel consumption from management perspective, which mainly include speed optimization, weather routing, efficient cargo operation, and trim optimization. As sailing speed is the most significant influencing factor, a considerable amount of literature has been focused on optimizing ship sailing speed to reduce fuel consumption, such as Fagerholt et al. (2010), Norstad et al. (2011), Yao et al. (2012), Wang and Meng (2012), Wang et al. (2013), Lee et al. (2015), Lindstad and Eskeland (2015), Song et al. (2015), Wang (2016) and Wang and Wang (2016). Weather routing helps ships to locally avoid rough sea and weather conditions in order to guarantee sailing safety as well as reduce fuel consumption. Studies on designing ship routes over a voyage based on weather information to realize fuel consumption reduction include Takashima et al. (2009), Shao et al. (2012), and Lin et al. (2013). IMO reported that trim optimization could reduce the main engine fuel consumption for most ship types by 0.5% to 3.0% (IMO, 2019). There is also research on developing trim optimization schemes for ship fuel consumption reduction, such as Reichel et al. (2014), Sherbaz and Duan (2014), Perera et al. (2015), and Moustafa et al. (2015). Proposing efficient ship cargo operation is often combined with fleet deployment and speed optimization, e.g., Xia et al. (2015) and Wang et al. (2015).

Over the period from 2016 through 2019, much more attention has been focused on developing two-phase optimization models for ship energy efficiency improvement. Generally, in the first phase, one or more models are developed for fuel consumption or weather conditions prediction under different situations; in the second phase, an optimization model is proposed for ship fuel consumption reduction over a voyage. Some typical two-phase models are presented as follows. Wang et al. (2016) proposed a real-time optimization model for a cruise ship which contained prediction of weather condition based on wavelet neural network (WNN) and determining the optimal engine speed based on the calculated ship resistance. Coraddu et al. (2017) developed a vessel trim optimization model for a tanker ship. The model included two parts: in the first part, a grey box model, which contained both mechanistic knowledge and historical data analysis, was proposed to predict the fuel consumption; in the second part, trim

optimization techniques were proposed. Lee et al. (2018) proposed a way to explore weather archive big data and optimize sailing speed for a container ship. First, the impact of weather conditions on ship fuel consumption was figured out by data mining methods. Then, speed optimization model was developed for the container. Du et al. (2019) presented a two-phase model for speed and trim optimization for a container. In the first phase, an ANN model was developed for estimating ship fuel consumption in different conditions. In the second phase, three countermeasures were put forward for reducing fuel consumption, including speed optimization, trim optimization as well as speed and trim optimization.

Although there are a growing number of studies on predicting and reducing ship fuel consumption, there are still considerable gaps existing in current literature. First, the literature has studied tankers, container ships, ferries, tugboats, and passenger ships. However, to the best of our knowledge, no model containing machine learning techniques for fuel consumption prediction and speed optimization are proposed for dry bulk ships. As the fuel prediction and optimization models are not universally applied (Lee et al., 2018; Banawan et al., 2013), it is necessary to develop such tailored model for a specific dry bulk carrier. Second, a large number of studies only include ship sailing speed as the input feature to predict ship fuel consumption. Actually, the determinants of fuel consumption are varied, including ship displacement and trim conditions as well as sea and weather conditions, but there is only a small number of studies considering these factors. Even if some studies take sea and weather information into account, the information is taken just from the noon report. Few studies have combined ship noon report with weather forecast, which could provide more comprehensive and accurate data. Third, most of the proposed machine learning models for ship fuel consumption prediction are based on ANNs. However, the development of ANN models usually requires a large number of training samples, and their structures are largely based on experience. In addition, tuning the parameters in ANNs can be difficult, and the prediction results are lack of interpretability. Moreover, the influencing degree of each input variable on the output variable is hard to figure out. Fourth, there are only a few pioneering studies on combining ship fuel prediction models and optimization models that can be put into practice to reduce fuel consumption and $CO_2$ emissions.

To bridge the gaps, we propose a two-stage model for a dry bulk ship based on ship noon report data and weather forecast data that contains (i) prediction of ship fuel consumption under different sailing speed, cargo, wind, swell, wind waves, and current conditions by adopting a random forest regressor, which is an ensemble learning method for regression based on multiple decision trees, and (ii) development of a speed optimization model to minimize ship fuel consumption over a voyage while

251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288

guaranteeing the estimated arrival time to the destination port based on the prediction results at the first stage. Compared with traditional statistical regression models, the advantages of using random forest regressors are that they are able to deal with high-dimensional data and make more accurate predictions. Compared to some other machine learning models, including ANNs, they are easier and faster to be implemented with more interpretable results and the influence degree of the features on the target variable can be generated, which can be used for feature selection.

**3. Data description**

**3.1 Ship noon report**

Noon report of a ship is a ship voyage report data sheet prepared by the ship's captain on a daily basis (usually at noon). Many attributes of the ship's sailing behavior are recorded, such as ship geographic location, distance travelled since last report, average propeller revolutions per minute (RPM), engine speed, sailing speed, total hold cargo and total deck cargo. In addition, sea and weather conditions of the recording time are also comprised, e.g., information on sea swell direction (coming direction of sea swell), sea swell height, sea current value (depth of sea current), sea current type (coming direction of sea current), wind force, wind direction (coming direction of wind), and sea temperature. It should be noted that although noon report data is the main source for ship fuel consumption and optimization research, the features contained in the report are limited and may vary among different reports. The factors used in other studies that also choose noon report as the data source for ship fuel consumption management are similar, such as "wind speed and direction, sea water temperature, air temperature, water depth, and wave height and direction" in the model calibrated by Pedersen and Larsen (2009), "displacement, wave direction and height, and wind force" in the model proposed by Meng et al. (2016), "displacement, wave direction, wind force and direction, sea current direction, sea water temperature, and trim" in the model developed by Du et al. (2019), "forward draft, aft draft, wind direction and wind Beaufort number" in the model presented by Yang et al. (2019).

**3.2 Description of ship voyage data**

The voyage data used in this study is the noon report data of a handy-size dry bulk ship with propeller diameter 5450mm, which was provided by an international shipping company. Time range of the voyage data is from 11[th] September 2017 to 27[th] February 2019. Initially, the voyage report data for the ship contains 738 data entries. To start with, we filter the data entries by choosing the records with ship conditions as "sailing at sea", "with cargo loaded" and sailing speed value no less than 5 knots. After preprocessing, there are 242 selected entries left in the entire case dataset. Then, we use the hourly fuel consumption of the ship as the target variable (which is calculated by
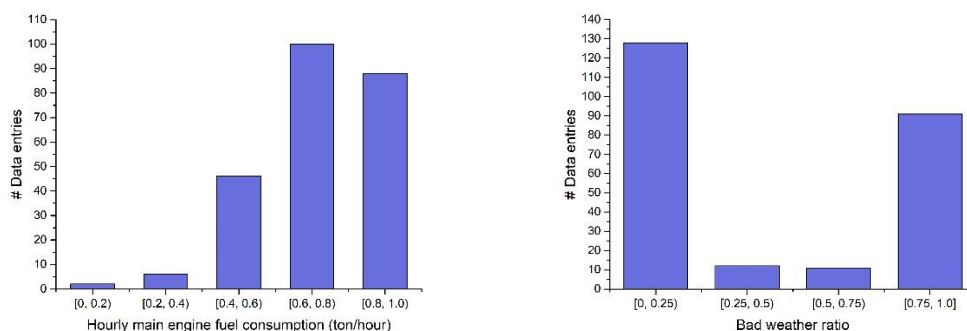
dividing the total fuel consumption by the total steaming hours) and delete the variables that are not suitable to be the input of the fuel consumption prediction model. Finally, 9 input variables are selected from the attributes in the voyage data, namely, bad weather ratio, ship sailing speed (knots), relative sea swell direction to ship's heading (°), sea swell height (m), sea current type, sea current value (m), relative wind direction to ship's heading (°), wind force (Beaufort force number), and total cargo weight (metric ton). Based on the recording time and location (longitude and latitude) provided by the noon report, we include two more attribute variables: height of combined wind waves and swell (m) and relative wind wave direction to ship's heading (°) downloaded from the European Centre for Medium-Range Weather Forecasts (ECMWF) (ECMWF, 2019). Eventually, the dataset contains 11 features as the input. Table 1 presents the statistical information of the variables in the case dataset. Figure 1 illustrates the distributions of the 242 data entries for the selected ship against the 11 input variables and the output variable.

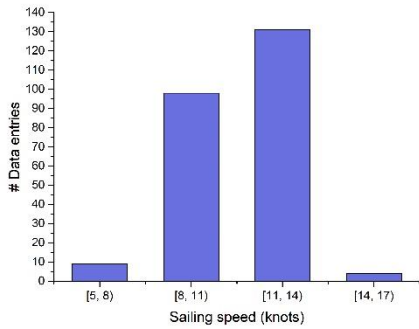Table 1. Description of the variables in the entire dataset

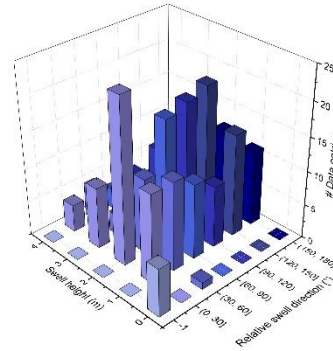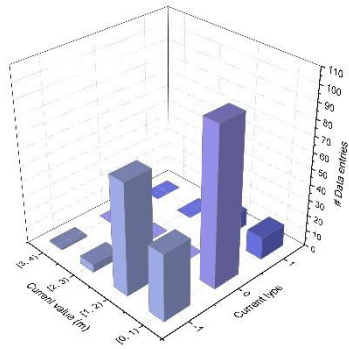| Variable name | Meaning | Unit | Max value | Min value | Mean value |
|---|---|---|---|---|---|
| Hourly main engine fuel consumption | Fuel consumption of main engine in an hour | MT/h | 0.9275 | 0.4139 | 0.7064 |
| Bad weather ratio | Steaming time in bad weather/total steaming time | \ | 1 | 0 | 0.4139 |
| Sailing speed | Ship average sailing speed | knots | 14.2 | 5.3 | 11.1021 |
| Relative sea swell direction | Direction of sea swell relative to ship's heading degree (-1 for no swell) | degree | 180 | 1 | 87.3440 |
| Sea swell height | Height of sea swell | meter | 4 | 0 | 1.9959 |
| Sea current type | Sea current against ship's heading (-1), with ship's heading (+1), no current (0) | \ | 1 | -1 | -0.3843 |
| Sea current value | Depth of sea current | meter | 4 | 0 | 0.4628 |
| Relative wind direction | Direction of wind relative to ship's heading degree (-1 for no wind) | degree | 180 | 0 | 99.8333 |
| Wind force | Measure of wind speed | Beaufort force number | 8 | 0 | 4.9008 |
| Total cargo weight | Sum of the weights of on-deck cargo and under-deck cargo | MT | 32741.5 | 11260 | 28685.06 |
| Combined wind waves and swell height | Height of the combination of wind waves and swell (-1 for no wave and swell) | meter | 9 | 0.1 | 2.1406 |
| Relative wind wave direction | Direction of wind wave relative to ship's heading degree (-1 for no wave) | degree | 179 | 0 | 81.1972 |

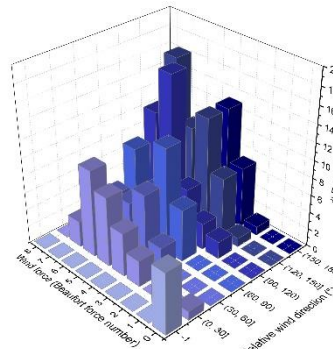(a) Distribution of hourly main engine fuel consumption
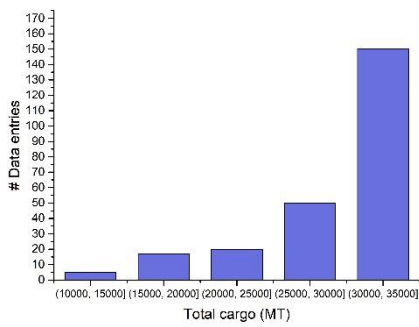
(b) Distribution of bad weather ratio
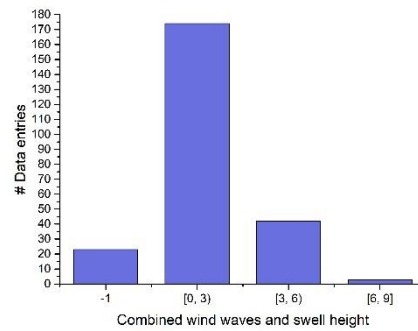


(c) Distribution of ship sailing speed

(d) Distribution of relative sea swell direction and sea swell height



(e) Distribution of sea current type and sea current value
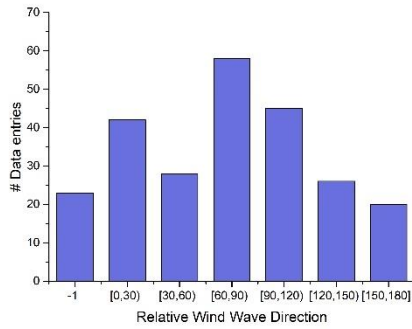
(f) Distribution of relative wind direction and wind force



(g) Distribution of total cargo weight

(h) Distribution of combined wind waves and swell height

(j) Distribution of relative wind wave
direction

Figure 1. Distribution of data entries in the entire dataset

In order to validate the performance of the proposed two-stage speed optimization model for ship fuel consumption reduction, we randomly select two 8-day continuous noon report data from the entire dataset which is used for numerical experiments. The noon report records from 16 Jan 2018 to 23 Jan 2018 and from 24 Dec 2018 to 31 Dec 2018 are selected, respectively. For the remaining 226 data entries, 80% of them are randomly selected to constitute the training set to develop the regression models, while the remaining 20% form the test set.

## 4. Development of tree-based models for ship fuel prediction

### 4.1 Introduction of Decision tree (DT) regression model

A decision tree (DT) is a supervised and tree-like decision support model which is widely used to predict both discrete valued output (*classification tree*) and continuous valued output (*regression tree*) (Myles et al., 2004). There are several nodes in a decision tree, and each node contains a certain number of input data entries. The output value of a node is the average output of all the comprised data entries. A decision tree consists of three types of nodes: root node (the topmost node), leaf node (which gives final prediction output), and internal node (node except for root and leaf node). The process of dividing a node into two successive nodes is called *splitting*. A feature and one of its corresponding value are chosen to split a node, and each splitting of a node requires finding out the *best split* based on some splitting criteria. In DT classifiers, common splitting criteria are *Gini impurity* and *information gain*. In DT regressor, common splitting criterion is *mean squared error* (MSE) (Friedman et al., 2001). A node being split is called parent node while the successive nodes are called child nodes.

Three widely used DT generation algorithms are ID3, C4.5, and CART (Classification and Regression Tree) (Loh, 2014). As both input and output data contain continuous valued data, we construct a DT regressor by adopting CART algorithm (Breiman et al., 1984). CART algorithm requires recursively and binarily splitting the

nodes, and a binary decision tree will be built. Originally, the construction process terminates when all the leaf nodes contain the data entries of the same output value. However, this usually means that the tree is extremely large and is heavily overfitted. To alleviate overfitting, termination criteria are preset to control tree dimension. Three commonly used criteria are presented as follows. It should be noted that the values of these decision tree parameters may vary from different training sets.

(a) The maximum depth of a tree (denoted by *max_depth*). The depth of a node in a decision tree is the number of nodes on a route from the root node to its parent node (the depth of root node is 0). The maximum depth of the tree is the maximum depth of all the nodes. A node cannot be further split if it reaches the maximum depth.

(b) The minimum number of data entries required to split a node (denoted by *min_samples_split*). If and only if a node contains data entries no less than *min_samples_split* can this node be further split.

(c) The minimum number of data entries required to be at a leaf node (denoted by *min_samples_leaf*). If and only if the number of data entries contained in both of the successive nodes split by the *best split* is no less than *min_samples_leaf* can the node be split.

Learning an optimal decision tree is known as an NP-complete problem (Laurent and Rivest, 1976; Naumov, 1991). Starting from splitting the root node, successive nodes are split in a depth-first manner until one of the termination criteria has been reached. Then, the next node for splitting is determined by retrospectively search for a node that can be further split. The algorithm terminates until there is no node that can be split. Main steps to generate a decision tree are described in Appendix A (Friedman et al., 2001; Harrington, 2012; Breiman, 2017).

**4.2 Introduction of random forest (RF) regression model**

Although the DT models are simple, intuitive, and interpretable, the main drawbacks of a single decision tree are that they are easy to get overfit (i.e., creating over-complex trees with poor generalization ability) and lack of robustness (i.e., small variations in the training data might result in a completely different tree being generated) (Ahmad et al., 2017). Ensemble learning is one of the popular ways to improve DT regressor performance. Ensemble methods contain multiple learning algorithms (called *weak learners*) and can obtain more desirable predictive performance than any of the constituent learning algorithms alone (Opitz and Maclin, 1999). There are two popular ensemble methods based on decision trees: *boosting* and *bagging*. In boosting, successive trees are dependent on the earlier trees, while in bagging, the trees are constructed using bootstrap sample of the training set (i.e., randomly selecting a certain number of samples from all the training samples with replacement) and the trees are independent on the other trees. Based on the bagging method, Breiman (2001) proposed

random forests by adding another layer of randomness: instead of considering all the data features to split the nodes in each DT included in the forest, a randomly generated subset of candidate features is used. Thus, apart from the abovementioned three parameters in DT regressor, there are two more parameters in RF regressor:

(d) The number of decision trees contained in the forest (denoted by *n_estimators*). Breiman (2001) proposed that adding more trees in the RF regressor will not suffer from overfitting. Instead, more trees have the ability to limit the value of generalization error.

(e) The number of features to consider when finding the *best split* of a node in each decision tree (denoted by *max_features*). The value of *max_features* should less than the total number of data features and the certain number of features are randomly selected at each splitting.

If CART based decision trees are the weak learners in a RF regressor, the main differences between constructing a DT regressor and a single decision tree in the RF regressor are twofold. (i) For a decision tree in RF regressor, bootstrap sampling from the entire training set to form a new training set is required; for a normal DT regressor, all the entries in the training set are used. (ii) For a decision tree in RF regressor, randomly selecting a subset of data features for splitting the nodes in each decision tree is required; for a normal DT regressor, all the features are considered when splitting each node. After a certain number of DTs are constructed, the RF regressor requires averaging the output values of all the tress as the prediction results (Liaw and Wiener, 2002). Compared with DT regressor, RF regressor has the advantages of robustness and lower variance (Siroky, 2009). For the detailed process of constructing an RF regressor, please refer to Breiman (2001), Biau and Scornet (2016), and Breiman (2017).

**4.3 Metrics for model validation**

In order to demonstrate the model performance in the test set, four typical regressor performance measures are adopted: *mean squared error* (MSE), *root mean squared error* (RMSE), *mean absolute error* (MAE), and *mean absolute percentage error* (MAPE). Denote the input variable vector by $x_e$, the predicted output value by $f(x_e)$, and the real output value by $y_e$. The total number of data entries in the test set is $N$. The definitions of MSE, RMSE, MAE, and MAPE are as follows:

$$MSE = \frac{1}{N} \sum_{e=1}^{N} [f(x_e) - y_e]^2 \tag{1}$$

$$RMSE = \sqrt{\frac{1}{N} \sum_{e=1}^{N} [f(x_e) - y_e]^2} \tag{2}$$

$$MAE = \frac{1}{N} \sum_{e=1}^{N} |f(x_e) - y_e| \tag{3}$$

$$MAPE = \frac{100\%}{N} \sum_{e=1}^{N} |\frac{f(x_e) - y_e}{y_e}|. \tag{4}$$

### 4.4 Construction and prediction results of DT and RF regression models

We adopt the scikit-learn machine learning library for Python to implement DT regressor and RF regressor based on CART algorithm (Pedregosa et al., 2011). The parameters for DT and RF regressors are set based on gird search method with five-fold cross validation as presented in Table 2. Except for those parameters, all the other parameters are set as the default values in scikit-learn library.

Table 2. Parameters used in DT and RF regressors

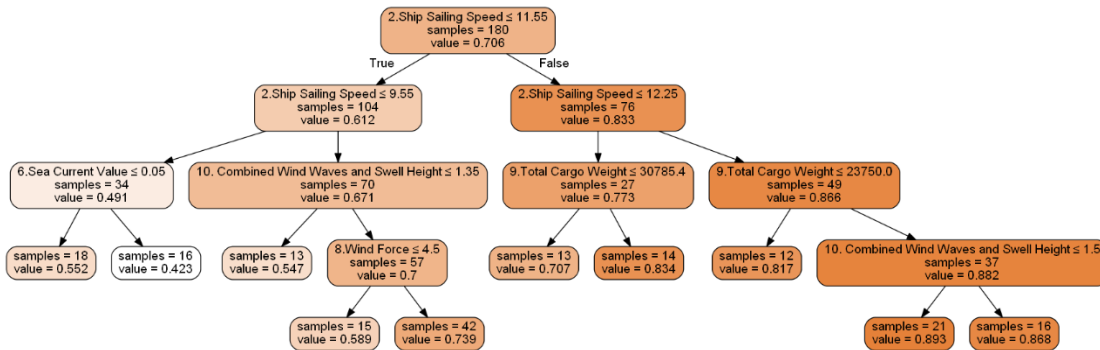| Parameter | Decision tree regressor | Random forest regressor |
|---|---|---|
| *max_depth* | 4 | 11 |
| *min_samples_split* | 5 | 2 |
| *min_samples_leaf* | 10 | 1 |
| *n_estimators* | / | 1000 |
| *max_features* | / | 4 |

The DT regressor model is visualized in Figure 2. For each root and internal node in the figure, the first row indicates the selected splitting variable and the splitting value. The second row shows the number of samples contained in the node. The third row is the output value of this node. For each leaf node, the first row is the number of samples contained in the node, and the second row is the final output value.



Figure 2. Visualization of the DT model for ship fuel prediction

The prediction performance of the two proposed regression models has also been compared with the popular machine learning based fuel consumption prediction methods in current literature. Three typical and popular regression models are selected for comparison: artificial neural network (ANN), least absolute shrinkage and selection operator (LASSO) regression, and support vector regression (SVR). ANN is a widely used machine learning model which contains a large number of highly interdependent processing elements called neurons. Usually, a typical ANN contains three layers of neurons: input layer, hidden layer, and output layer. LASSO is a linear regression

analysis method that can perform variable selection and regularization in order to improve regression performance. SVR is an application of support vector machine (SVM) to regression problems. The datasets used for training and testing the models are the same as those are used for the DT and RF regressors. It is worth mentioning that the LASSO and SVR models are implemented by adopting scikit-learn machine learning library for Python, and the parameters for the two models are tuned by gird search method with five-fold cross validation. The construction of the ANN model is similar to Du et al. (2019): five ANN models are constructed in MATLAB R2017a and the average of the outputs of the five ANN models is the prediction output. The prediction performance of the DT regressor, RF regressor, ANN, LASSO, and SVR models are on a daily basis, i.e., hourly fuel consumption data has been converted to fuel consumption for a day by considering steaming hours, are shown in Table 3. It can be seen that both of the tree-based regressors perform well on our test set and the RF regressor performs the best. Moreover, the RF regressor outperforms the DT regressor regarding every metric.

Table 3. Performance of the five regression models on test set

| Model/Metric | MSE | RMSE | MAE | MAPE |
|---|---|---|---|---|
| DT regressor | 6.16 | 2.48 | 1.74 | 11.33% |
| **RF regressor** | **3.17** | **1.78** | **1.21** | **7.91%** |
| ANN | 5.67 | 2.38 | 1.68 | 11.95% |
| LASSO | 5.60 | 2.37 | 1.72 | 11.51% |
| SVR | 9.37 | 3.06 | 2.55 | 15.47% |

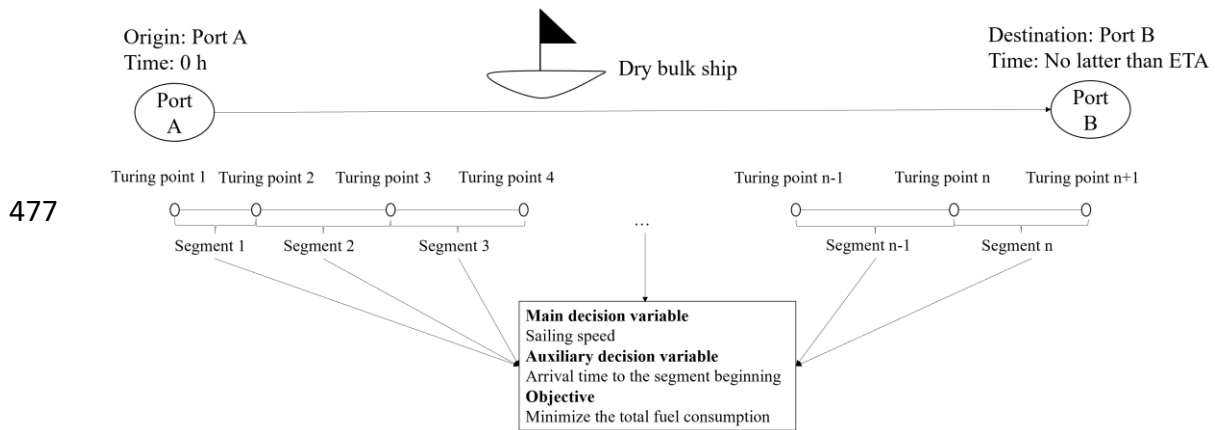## 5. Development of ship speed optimization model

### 5.1 Problem description

As mentioned in the introduction part, shipping companies have a strong motivation to carefully plan ship sailing speed during a voyage to reduce bunker consumption and comply with the environmental protection conventions. Based on the proposed RF regressor, which is able to predict fuel consumption of the dry bulk ship under different sailing speeds, total cargo weight and sea and weather conditions, we develop a speed optimization model between two ports while guaranteeing ship estimated time of arrival (ETA) to the destination port.

To develop the sailing speed optimization model, we consider a situation when this dry bulk ship sails from an origin port A to a destination port B along a fixed path which the captain is quite familiar with. The loaded cargo of this ship is pre-determined and fixed during the voyage and the sea and weather conditions can be obtained via forecasts 5 to 7 days in advance. Due to the dynamic conditions at sea, the whole path can be divided into several segments, and in each segment, we assume that the

international conventions that the ship needs to obey as well as sea and weather conditions can be viewed as identical. A lower bound and an upper bound of permitted ship sailing speed are also the same in one segment. The range of the allowable sailing speed is determined by many factors, especially the sea and weather conditions in the segment. For example, Tsou and Cheng (2013) adopted a formula to calculate a ship's allowable maximum speed while navigation in storm conditions based on wave height and wave direction to ensure navigational safety. The ship departure time from port A is 0, and the ETA of port B is no later than the latest allowable arrival time. Two questions need to be addressed: when to adjust the sailing speed (referred to as the time of speed turning point) and what speed should be adjusted to (referred to as adjusted speed). The objective of the model is to minimize the total fuel consumption of this dry bulk ship over the whole voyage by determining the sailing speed in each segment. Except for sailing speed, external factors that influence ship fuel consumption are the same in one segment and thus the optimal speed should be the same in a segment. Thus, it can be justified that the speed turning points can only occur at the beginning of a segment. An illustration of the optimization problem is presented in Figure 3.



Figure 3. An illustration of the problem

**5.2 Development of a mathematical model**

The notation of the mathematical model is defined as follows.

| Sets and indices | |
| --- | --- |
| $n$ | Total number of path segments |
| $i$ | Index of a path segment, $i \in \{1,...,n+1\}$. Segment $n+1$ presents the end of segment $n$, i.e., port B |
| $I$ | Set of all path segments, $I = \{1,...,n\}$ |

| Parameters | |
| --- | --- |
| $c_i$ | Ship total loaded cargo and sea and weather conditions in segment $i$ |
| $v_0$ | Ship speed before departure |
| $v_{i,c_i}^{max}$ | Maximum allowable speed when sailing in segment $i$ with the loaded |

16

| | |
|---|---|
| | cargo and sea and weather conditions as $c_i$ (knots) |
| $v_{i,c_i}^{min}$ | Minimum allowable speed when sailing in segment $i$ with the loaded cargo and sea and weather conditions as $c_i$ (knots) |
| $f^{RF}(v,c)$ | Predicted ship fuel consumption (ton/hour) by using the proposed RF model when sailing speed is $v$ and ship total loaded cargo and sea and weather conditions are $c$ |
| $L_i$ | Path length of segment $i$ (nm) |
| $T_{max}$ | Latest allowable arrival time to the destination port |

482

| | |
|---|---|
| **Main decision variables** | |
| $v_i$ | Ship sailing speed in segment $i$ (knots) |
| **Auxiliary decision variable** | |
| $t_i$ | Arrival time to the beginning of segment $i$, $t_1 = 0$. $t_{n+1}$ is the arrival time to the end of segment $n$, i.e., the arrival time of port B. |

483
484 The speed optimization problem can be formulated by using Model **M1** based on the parameters and decision variables.

485 [**M1**]

486
$$\min \sum_{i=1}^{n}(f^{RF}(v_i,c_i) \times \frac{L_i}{v_i}) \tag{5}$$

487 subject to:

488
$$t_{i+1} = t_i + \frac{L_i}{v_i}, \forall i \in I \tag{6}$$

489
$$t_{n+1} \leq T_{max} \tag{7}$$

490
$$v_0 = 0 \tag{8}$$

491
$$v_{i,c_i}^{min} \leq v_i \leq v_{i,c_i}^{max}, \forall i \in I \tag{9}$$

492
$$t_i \geq 0, \forall i \in I \bigcup \{n+1\}. \tag{10}$$

493 Objective (5) minimizes ship fuel consumption over the voyage. Constraint (6)
494 indicates the relationship between the arrival time to the beginning of the previous
495 segment and that of the next segment. Constraint (7) ensures the ship arrival time to the
496 destination port is no later than the allowable arrival time. Constraint (8) ensures the
497 sailing speed before departure is 0. Constraint (9) guarantees the lower and upper
498 bounds of the sailing speed in each segment. Constraint (10) grantees that the arrival
499 time to the beginning of every segment is nonnegative. M1 cannot be solved directly
500 by the off-the-shelf optimizers, thus we linearize the model in the next section.

## 5.3 Linearization of model M1

Given the maximum and minimum allowable sailing speeds, we can discretize the speed values with 0.1 knot as an interval. Specifically, given $v_{i,c_i}^{\max}$ and $v_{i,c_i}^{\min}$ in segment $i$ respectively, as we discretize the sailing speeds with 0.1 as an interval, we have the sailing speed parameters $v_i^1 = v_{i,c_i}^{\min}$, $v_i^2 = v_i^1 + 0.1$,..., $v_i^{u_i} = v_{i,c_i}^{\max}$ and a specific sailing speed as $v_i^u = \{v_i^1, v_i^2, ..., v_i^{u_i}\}$. We further introduce a binary decision variable $y_i^u \in \{0,1\}$, and if $v_i = v_i^u$, $y_i^u = 1$; otherwise $y_i^u = 0$. The new main decision variable is $y_i^u$, and the auxiliary decision variable is $t_i$. Based on the new parameters and decision variables, we can convert model M1 to model M2.

[**M2**]

$$\min \sum_{i=1}^{n} \sum_{u=1}^{u_i} [y_i^u \times (f^{RF}(v_i^u, c_i) \times \frac{L_i}{v_i^u})] \tag{11}$$

subject to:

$$t_{i+1} = t_i + \sum_{u=1}^{u_i} (\frac{L_i}{v_i^u} \times y_i^u), \forall i \in I \tag{12}$$

$$t_{n+1} \leq T_{\max} \tag{13}$$

$$v_0 = 0 \tag{14}$$

$$y_i^u \in \{0,1\}, \forall i \in I, u \in \{1,2,...,u_i\} \tag{15}$$

$$\sum_{u=1}^{u_i} y_i^u = 1, \forall i \in I \tag{16}$$

$$t_i \geq 0, \forall i \in I \bigcup \{n+1\} \tag{17}$$

Model M2 is equivalent to M1 and is a mixed-integer linear programming (MIP) model, which can be solved by the off-the-shelf optimizers, such as CPLEX.

## 6. Computational experiments

### 6.1 Prediction of ship fuel consumption

In this section, we adopt the RF regressor developed in Section 4 to predict fuel consumption of the dry bulk ship during two continuous 8-day sailing voyages (denoted by voyage 1 and voyage 2, respectively). The total sailing distance of voyage 1 is 2001.2 nautical miles with the total sailing time as 195 hours, and the total fuel consumption is 126.36 tons. The total sailing distance of voyage 2 is 1946.4 nautical miles with the total sailing time as 192 hours, and the total fuel consumption is 114.92 tons. The total cargo weight as well as sea and weather conditions in each sailing segment of each voyage are presented in Table 4 and Table 5.

Table 4. Initial ship sailing information of voyage 1

| Day | Steaming hour (h) | Total fuel (tons) | Speed (knots) | Hourly fuel consumption (tons/h) | Bad weather ratio | Relative sea swell direction (degree) | Sea swell height (m) | Sea current type | Sea current value | Wind force (Beaufort force number) | Relative wind direction (degree) | Total cargo weight (MT) | Combined Wind Waves and Swell Height (m) | Relative Wind Wave Direction (degree) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 24 | 15.69 | 10 | 0.65375 | 1 | 160 | 2 | 0 | 0 | 6 | 160 | 26605 | 1.7 | 136 |
| 2 | 24 | 15.42 | 9.9 | 0.64250 | 1 | 89 | 2 | 0 | 0 | 6 | 111.5 | 26605 | 3.0 | 92.3 |
| 3 | 25 | 16.22 | 9.8 | 0.64880 | 0.40 | 44 | 1 | 0 | 0 | 4 | 89 | 26605 | 1.2 | 6.5 |
| 4 | 24 | 15.43 | 10.2 | 0.64292 | 0.25 | 82 | 2 | 0 | 0 | 4 | 8 | 26605 | 2.4 | 52.8 |
| 5 | 25 | 16.10 | 10.4 | 0.64400 | 0.80 | 170 | 2 | 0 | 0 | 5 | 170 | 26605 | 0.8 | 170.3 |
| 6 | 24 | 15.45 | 10.3 | 0.64375 | 1 | 106 | 2 | 0 | 0 | 6 | 128.5 | 26605 | 3.3 | 80.5 |
| 7 | 24 | 15.54 | 10.9 | 0.64750 | 0.50 | 89 | 2 | 0 | 0 | 5 | 91 | 26605 | 1.2 | 28.1 |
| 8 | 25 | 16.51 | 10.6 | 0.66040 | 1 | 99 | 2 | 0 | 0 | 7 | 144 | 26605 | 2.3 | 90.3 |

Table 5. Initial ship sailing information of voyage 2

| Day | Steaming hour (h) | Total fuel (tons) | Speed (knots) | Hourly fuel consumption (tons/h) | Bad weather ratio | Relative sea swell direction (degree) | Sea swell height (m) | Sea current type | Sea current value | Wind force (Beaufort force number) | Relative wind direction (degree) | Total cargo weight (MT) | Combined Wind Waves and Swell Height (m) | Relative Wind Wave Direction (degree) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 24 | 14.45 | 11.2 | 0.60208 | 0 | 113 | 1 | -1 | 0.6 | 3 | 90.5 | 31948 | 1.9 | 93.0 |
| 2 | 24 | 14.41 | 10.8 | 0.60042 | 0 | 87 | 1 | -1 | 0.6 | 4 | 154.5 | 31948 | 2.0 | 76.9 |
| 3 | 24 | 14.35 | 9.2 | 0.59792 | 0.75 | 132 | 2 | -1 | 1 | 6 | 177 | 31948 | 1.9 | 106.4 |
| 4 | 24 | 14.36 | 8.9 | 0.59833 | 1 | 177 | 2 | -1 | 1.5 | 6 | 160.5 | 31948 | 1.6 | 121.6 |
| 5 | 24 | 14.33 | 9.7 | 0.59708 | 0.17 | 95 | 2 | -1 | 1.5 | 4 | 117.5 | 31948 | 1.5 | 111.8 |
| 6 | 24 | 14.36 | 10.5 | 0.59833 | 0 | 100 | 2 | -1 | 1 | 4 | 100 | 31948 | 1.8 | 104.6 |
| 7 | 24 | 14.45 | 11.5 | 0.60208 | 0.04 | 145 | 1 | 0 | 0 | 5 | 100 | 31948 | 2.3 | 106.1 |
| 8 | 24 | 14.21 | 9.3 | 0.59208 | 1 | 145 | 2 | -1 | 1.5 | 6 | 100 | 31948 | 3.0 | 113.0 |

We use the total 226 data entries (i.e., all the data entries in the entire dataset except for the 16 records used to validate the optimization model) to construct the RF regressor. Due to the lack of extreme valued data, predicting fuel consumption under too large or too small speed values is highly likely to suffer from inaccuracy (Freidman et al., 2001). Thus, we make predictions on fuel consumption with speed values ranging between 10% and 90% from small to large in the training set, i.e., we exclude the 10% smallest and 10% largest speed values. The selected speed values are from 8.9 to 13.3 knots. The fuel consumption prediction results are presented in Figure 4. For the 16 validation records, we only have the real output value under the given speed. The performance of the RF regressor on predicting the fuel consumption under the given speed of the 16 records are given in Table 6. It can be seen that the predicted fuel consumption under the given speed is higher than the real fuel consumption.
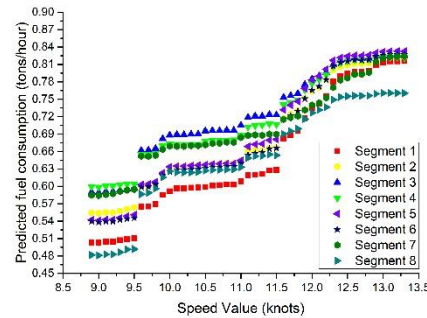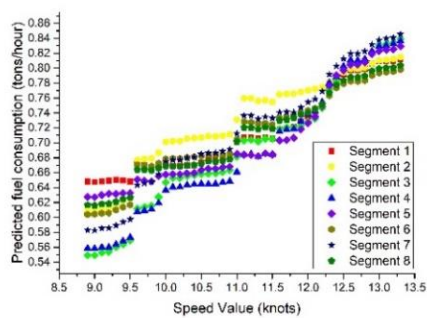
(a) Prediction results of voyage 1      (b) Prediction results of voyage 2

Figure 4. Fuel consumption prediction results of the voyages

Table 6. RF regressor performance on predicting fuel
consumption of the two voyages

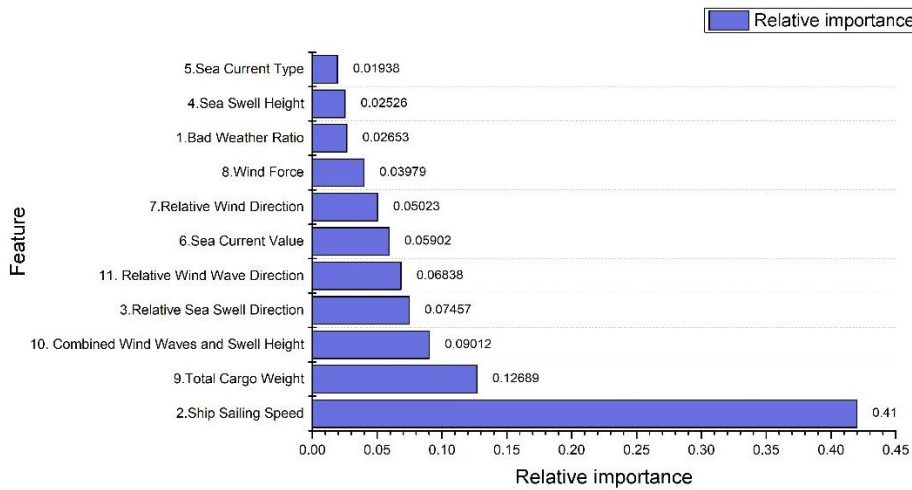| Path NO | Real total fuel consumption | Predicted total fuel consumption | MSE | RMSE | MAE | MAPE |
|---|---|---|---|---|---|---|
| 1 | 126.36 | 129.6335 | 0.5186 | 0.7201 | 0.6285 | 4.00% |
| 2 | 114.92 | 116.6830 | 1.5702 | 1.2531 | 0.8950 | 6.24% |

Apart from predicting fuel consumption under different speed values in each segment, the RF regressor is also able to illustrate the feature importance of the input variables when predicting fuel consumption. The feature importance generated by the 230 data entries is shown in Figure 5.



Figure 5. Relative importance of the input features

Figure 5 indicates that ship sailing speed is the most significant influencing factor on ship fuel consumption, which allies with most of the current literature. Total cargo weight also has a great influence on fuel consumption. Regarding sea and weather conditions, combined wind wave and swell height, relative sea swell, and relative wind wave direction can have more impact, while sea current type has the least impact on ship fuel consumption.

**6.2 Validation of speed optimization model**

We validate the performance of the proposed mathematical model by adopting CPLEX optimizer to find the optimal sailing speed in each segment which can minimize the total fuel consumption. A laptop (Intel Core i7, 2.20GHz, Memory 16.0G) is used to conduct the experiment with the programming language C++. The selected sailing

speed, sailing time, and fuel consumption in each segment of the two voyages are shown in Table 7 and Table 8.

Table 7. Ship sailing information in voyage 1 after optimization

| Day | Sailing hours (h) | Sailing speed (knots) | Hourly fuel consumption (tons/h) | Total fuel consumptions (tons) |
|-----|-------------------|-----------------------|----------------------------------|-------------------------------|
| 1 | 18.0451 | 13.3 | 0.812160 | 14.6555 |
| 2 | 17.8647 | 13.3 | 0.814638 | 14.5532 |
| 3 | 25.7895 | 9.5 | 0.569093 | 14.6766 |
| 4 | 22.4587 | 10.9 | 0.648483 | 14.5641 |
| 5 | 22.6087 | 11.5 | 0.684186 | 15.4685 |
| 6 | 18.5865 | 13.3 | 0.798043 | 14.8328 |
| 7 | 21.6198 | 12.1 | 0.758690 | 16.4027 |
| 8 | 19.9248 | 13.3 | 0.803543 | 16.0104 |
| Total | 166.898 | / | / | 121.1638 (-4.11%) |

Table 8. Ship sailing information in voyage 2 after optimization

| Day | Sailing hours (h) | Sailing speed (knots) | Hourly fuel consumption (tons/h) | Total fuel consumptions (tons) |
|-----|-------------------|-----------------------|----------------------------------|-------------------------------|
| 1 | 28.2947 | 9.5 | 0.510766 | 14.4520 |
| 2 | 22.5391 | 11.5 | 0.666378 | 15.0196 |
| 3 | 16.6015 | 13.3 | 0.825215 | 13.6998 |
| 4 | 18.5739 | 11.5 | 0.706721 | 13.1266 |
| 5 | 24.5053 | 9.5 | 0.550893 | 13.4998 |
| 6 | 26.5263 | 9.5 | 0.545664 | 14.4745 |
| 7 | 24 | 11.5 | 0.689559 | 16.5494 |
| 8 | 23.4947 | 9.5 | 0.491945 | 11.5581 |
| Total | 184.536 | / | / | 112.380 (-2.21%) |

It takes only 0.05 and 0.03 second respectively to find the optimal solutions for case 1 and case 2. After speed optimization, the dry bulk ship consumes 121.1638 tons of fuel to complete voyage 1 and 112.38 tons of fuel to complete voyage 2 while guaranteeing the arrival time to the destination. Compared with the real fuel consumption of 126.36 and 114.92 tons in voyage 1 and 2, the ship can save 4.11% and 2.21% of total fuel consumption after speed optimization. It should be noted that the predicted fuel consumption is a little higher than the real fuel consumption. If we compare the fuel consumption after speed optimization and the predicted total fuel consumption in these two voyages (i.e. 129.6335 and 116.6830), we can conclude that 6.53% and 3.69% of fuel consumption can be reduced, respectively. We can then conclude that the two-stage fuel consumption prediction and speed optimization model can help the bulk carrier ship to save 2% to 7% fuel to complete an 8-day voyage. Note that as the data for training the RF regressor is limited and there can be inaccuracy in fuel consumption prediction, the savings in fuel consumption may have variations.

**7. Extension and future research**

Although numerous studies on ship fuel management are conducted based on ship noon report, the data of ship sailing information and sea and weather conditions provided by the noon report is actually limited (as discussed in Section 3.1) and the time resolution is low: usually only one record for 24 hours. In order to make precise fuel consumption prediction, which is the foundation of efficient ship fuel management, one possible way is to incorporate more data features from other data sources, such as sea and weather data from weather forecast website. For example, water depth can be included for considering the influence of shallow water on fuel consumption. Besides, temperature and salinity of water can also have an effect. Another possible way is to combine noon report with ship sensory data, which can provide more ship sailing features such as trim and draft condition with much lower time interval. Combining sensory data can also help to develop sailing speed optimization model. If more data can be obtained for a day, the division of sailing segment can be more flexible, e.g., by the length of a fixed time such as 3h, 6h, or 12h voyage at the calm water set speed. In addition, if more ship sailing information is accessible, it is easier to combine sea and weather data with ship sailing data. In the current noon report, the sailing distance of one day is usually more than 200 nm. As only one record is generated for each day, the sea and weather conditions are viewed as identical in the whole sailing distance covered by a whole day. However, sea and weather data can usually have a given resolution. For example, if the resolution of sea and weather forecasting data is 0.5°, the associated arc length is 30 nm. In addition, the weather forecast usually renews every few hours. For example, the weather forecast provided by ECMWF renews every 6 hours. If we can have 4–6 reports each day, the ship sailing data and sea and weather data can be combined. Therefore, accurate and practical fuel consumption prediction model and sailing speed optimization model can be proposed.

## 8. Conclusion

Shipping companies are developing stronger motivations to improve ship fuel energy efficiency for the purpose of complying with environmental conventions and increasing their profits. This study proposes a two-stage fuel consumption prediction and reduction model for a dry bulk ship to improve its energy efficiency based on the noon report data. More specifically, at the first stage, a random forest regression model is developed to predict the dry bulk ship's fuel consumption under different total carried cargo, sea, and weather conditions. It is also validated that the proposed RF regressor outperforms the widely used machine learning models such as ANN, SVR, and LASSO for ship fuel consumption prediction. At the second stage, a speed optimization model is proposed based on the fuel consumption prediction results at the first stage. The objective of the optimization model is to minimize the total fuel consumption of the dry

bulk ship over a voyage which contains several segments by deciding the sailing speed in each segment. The model is a mixed integer programming model which can be efficiently solved by CPLEX. In the computational experiments, we use two 8-day sailing voyage reports to test the performance of the two-stage model. The results show that the proposed model can save fuel consumption to 2%–7% compared with the real situation, which can also lead to significant $CO_2$ emissions reduction. In addition, the influence degree of the input features on the total fuel consumption is also generated. Similar to other related studies, it is indicated that ship sailing speed is the dominant factor of ship fuel consumption, then followed by total carried cargo. Regarding sea and weather conditions, combined wind wave and swell height, relative sea swell and wind wave directions can also have remarkable impact, while the current type has the least influence on ship fuel consumption. This paper considers the relationship between ship sailing speed and fuel consumption rate in a non-analytical form, which improves the common understanding about fuel consumption management. The proposed model is one of the pioneering models which combine a machine learning model with an optimization model in ship fuel consumption prediction and reduction. Based on the model, shipping companies are able to finer plan the daily sailing speed of their ships in order to reduce fuel consumption and $CO_2$ emissions.

**Acknowledgment**

**Reference**

Ahmad, M. W., Mourshed, M., Rezgui, Y., 2017. Trees vs neurons: comparison between random forest and ANN for high–resolution prediction of building energy consumption. Energy and Buildings 147, 77–89.

Alpaydin, E., 2009. Introduction to Machine Learning. MIT Press, Cambridge.

Andersen, P., Borrod, A. S., Blanchot, H., 2005. Evaluation of the service performance of ships. Marine Technology 42(4), 177–183.

Banawan, A. A., Mosleh, M., Seddiek, I. S., 2013. Prediction of the fuel saving and emissions reduction by decreasing speed of a catamaran. Journal of Marine Engineering & Technology 12(3), 40–48.

Beşikçi, E. B., Arslan, O., Turan, O., Ölçer, A. I., 2016. An artificial neural network based decision support system for energy efficient ship operations. Computers & Operations Research 66, 393–401.

Bialystocki, N., Konovessis, D., 2016. On the estimation of ship's fuel consumption and speed curve: a statistical approach. Journal of Ocean Engineering and Science 1(2), 157–166.

Biau, G., Scornet, E., 2016. A random forest guided tour. Test 25(2), 197–227.

Bishop, C. M., 2006. Pattern Recognition and Machine Learning. Springer Publisher, Berlin.

Bocchetti, D., Lepore, A., Palumbo, B., Vitiello, L., 2013. A statistical control of the ship fuel consumption. In: International Conference on the Design, Construction and Operation of Passenger Ships Proceedings, 91–96.

Bocchetti, D., Lepore, A., Palumbo, B., Vitiello, L., 2015. A statistical approach to ship fuel consumption monitoring. Journal of Ship Research 59(3), 162–171.

Breiman, L., 2001. Random forests. Machine Learning 45(1), 5–32.

Breiman, L., 2017. Classification and Regression Trees. Routledge Publisher, London.

Breiman, L., Friedman, J., Stone C. J., Olshen R. A., 1984. Classification and Regression Trees. Taylor & Francis, Abingdon.

Coraddu, A., Oneto, L., Baldi, F., Anguita, D., 2015. Ship efficiency forecast based on sensors data collection: improving numerical models through data analytics. In: IEEE Conference on OCEANS Proceedings,1–10.

Coraddu, A., Oneto, L., Baldi, F., Anguita, D., 2017. Vessels fuel consumption forecast and trim optimisation: a data analytics perspective. Ocean Engineering 130, 351–370.

Corbett, J. J., Wang, H., Winebrake, J. J., 2009. The effectiveness and costs of speed reductions on emissions from international shipping. Transportation Research Part D: Transport and Environment 14(8), 593–598.

Du, Y., Meng, Q., Wang, S., Kuang, H., 2019. Two-phase optimal solutions for ship

702       speed and trim optimization over a voyage using voyage report data. Transportation
703       Research Part B: Methodological 122, 88–114.

704   ECMWF, 2019. <https://apps.ecmwf.int/datasets/data/interim-full-daily/levtype=sfc/>
705       (Accessed 8 Oct 2019).

706   Erto, P., Lepore, A., Palumbo, B., Vitiello, L., 2015. A procedure for predicting and
707       controlling the ship fuel consumption: its implementation and test. Quality and
708       Reliability Engineering International 31(7), 1177–1184.

709   Fagerholt, K., Laporte, G., Norstad, I., 2010. Reducing fuel emissions by optimizing
710       speed on shipping routes. Journal of the Operational Research Society 61(3), 523–
711       529.

712   Friedman, J., Hastie, T., Tibshirani, R., 2001. The Elements of Statistical Learning.
713       Springer Publisher, Berlin.

714   Goldstein, H., 2011. Multilevel Statistical Models. John Wiley & Sons, New York.

715   Haranen, M., Pakkanen, P., Kariranta, R., Salo, J., 2016. White, grey and black-box
716       modelling in ship performance evaluation. In Proceedings of the 1st Hull
717       Performance & Insight Conference, 115-127.

718   Harrington, P., 2012. Machine Learning in Action. Manning Publisher, New York.

719   Hasselaar, T. W. F., 2011. An investigation into the development of an advanced ship
720       performance monitoring and analysis system. Doctoral dissertation, Newcastle
721       University.

722   Holtrop, J., 1977. Statistical analysis of performance test results. International
723       Shipbuilding Progress 24, 23-28.

724   Holtrop, J., 1978. Statistical data for the extrapolation of model performance
725       tests. International Shipbuilding Progress 25, 1-5.

726   IMO, 2011. Amendments to the annex of the protocol of 1997 to amend the
727       International Convention for the Prevention of Pollution from Ships, 1973, as
728       modified by the protocol of 1978 relating thereto.
729       <http://www.imo.org/en/KnowledgeCentre/IndexofIMOResolutions/Marine
730       Environment-Protection-Committee-(MEPC)/Documents/MEPC.203(62).pdf>.
731       Accessed (28.6.2019).

732   IMO, 2014. Third IMO Greenhouse Gas Study 2014. <http://www.imo.org/en/A
733       bout/Pages/IMODocuments.aspx>. Accessed (10.7.2019)

734   IMO, 2019. Trim and draft optimization. <https://glomeep.imo.org/technology/trim–
735       and–draft–optimization/>. (Accessed 3.7.2019).

736   Journée, J.M.J., Rijke, R.J., Verleg, G.J.H., 1987. Marine performance surveillance with
737       a personal computer (Technical Report No. 753- P). Delft University of Technology,
738       Ship Hydromechanics Laboratory, Delft, The Netherlands.

739   Kee, K. K., Simon, L., Renco, Y., 2018. Prediction of ship fuel consumption and speed

curve by using statistical method. Journal of Computer Science & Computational Mathematics 8(2), 19–24.

Kristensen, H. O., Lützen, M., 2012. Prediction of resistance and propulsion power of ships. Clean Shipping Currents 1(6), 1-52.

Kwon, Y. J., 1981. The effect of weather, particularly short sea waves, on ship speed performance (Doctoral dissertation, Newcastle University).

Laurent, H., Rivest, R. L., 1976. Constructing optimal binary decision trees is NP-complete. Information Processing Letters 5(1), 15–17.

Lee, C. Y., Lee, H. L., Zhang, J., 2015. The impact of slow ocean steaming on delivery reliability and fuel consumption. Transportation Research Part E: Logistics and Transportation Review 76, 176–190.

Lee, H., Aydin, N., Choi, Y., Lekhavat, S., Irani, Z., 2018. A decision support system for vessel speed decision in maritime logistics using weather archive big data. Computers & Operations Research 98, 330–342.

Leifsson, L. Þ., Sævarsdóttir, H., Sigurðsson, S. Þ., Vésteinsson, A., 2008. Grey-box modeling of an ocean vessel for operational optimization. Simulation Modelling Practice and Theory 16(8), 923–932.

Liaw, A., Wiener, M., 2002. Classification and regression by random forest. R News 2(3), 18–22.

Lin, Y. H., Fang, M. C., Yeung, R. W., 2013. The optimization of ship weather-routing algorithm based on the composite influence of multi-dynamic elements. Applied Ocean Research 43, 184–194.

Lindstad, H., Eskeland, G. S., 2015. Low carbon maritime transport: how speed, size and slenderness amounts to substantial capital energy substitution. Transportation Research Part D: Transport and Environment 41, 244–256.

Lo, H. K., McCord, M. R., 1995. Routing through dynamic ocean currents: general heuristics and empirical results in the gulf stream region. Transportation Research Part B: Methodological 29(2), 109-124.

Loh, W. Y., 2014. Fifty years of classification and regression trees. International Statistical Review 82(3), 329–348.

Lu, R., Turan, O., Boulougouris, E., 2013. Voyage optimization: prediction of ship specific fuel consumption for energy efficient shipping. In: Conference on Low Carbon Shipping Proceedings, 1–11.

Lu, R., Turan, O., Boulougouris, E., Banks, C., Incecik, A., 2015. A semi-empirical ship operational performance prediction model for voyage optimization towards energy efficient shipping. Ocean Engineering 110, 18-28.

MAN Diesel & Turbo, 2011. Basic Principles of Ship Propulsion. MAN Diesel & Turbo, Copenhagen, Denmark.

778    Meng, Q., Du, Y., Wang, Y., 2016. Shipping log data based container ship fuel efficiency
779        modeling. Transportation Research Part B: Methodological 83, 207–229.

780    Moustafa, M. M., Yehia, W., Hussein, A. W., 2015. Energy efficient operation of bulk
781        carriers by trim optimization. In: International Conference on Ships and Shipping
782        Research Proceedings, 484–493.

783    Myles, A. J., Feudale, R. N., Liu, Y., Woody, N. A., Brown, S. D., 2004. An introduction
784        to decision tree modeling. Journal of Chemometrics 18(6), 275–285.

785    Naumov, G. E., 1991. NP-completeness of problems of construction of optimal decision
786        trees. Soviet Physics 1, 270–271.

787    Neter, J., Kutner, M. H., Nachtsheim, C. J., Wasserman, W., 1996. Applied Linear
788        Statistical Models. McGraw–Hill Irwin Publisher, Chicago.

789    Norstad, I., Fagerholt, K., Laporte, G., 2011. Tramp ship routing and scheduling with
790        speed optimization. Transportation Research Part C: Emerging Technologies 19(5),
791        853–865.

792    Opitz, D., Maclin, R., 1999. Popular ensemble methods: an empirical study. Journal of
793        Artificial Intelligence Research 11, 169–198.

794    Pedersen, B. P., Larsen, J., 2009. Prediction of full–scale propulsion power using
795        artificial neural networks. In: the 8th International Conference on Computer and IT
796        Applications in the Maritime Industries (COMPIT'09) Proceedings, 537–550.

797    Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel
798        M., Prettenhofer P., Weiss R., Dubourg V., Passos A., Cournapeau D., Brucher M.,
799        Perrot M., Duchesnay E., 2011. Scikit-learn: Machine learning in Python. Journal
800        of Machine Learning Research 12, 2825–2830.

801    Perera, L. P., Mo, B., Kristjánsson, L. A., 2015. Identification of optimal trim
802        configurations to improve energy efficiency in ships. In: the 10th IFAC Conference
803        on Manoeuvring and Control of Marine Craft Proceedings, 267–272.

804    Petersen, J. P., Jacobsen, D. J., Winther, O., 2012a. Statistical modelling for ship
805        propulsion efficiency. Journal of Marine Science and Technology 17(1), 30–39.

806    Petersen, J. P., Winther, O., Jacobsen, D. J., 2012b. A machine-learning approach to
807        predict main energy consumption under realistic operational conditions. Ship
808        Technology Research 59(1), 64–72.

809    Petursson, S., 2009. Predicting optimal trim configuration of marine vessels with
810        respect to fuel usage. Doctoral dissertation, University of Iceland.

811    Psaraftis, H. N., Kontovas, C. A., 2013. Speed models for energy-efficient maritime
812        transportation: a taxonomy and survey. Transportation Research Part C: Emerging
813        Technologies 26, 331–351.

814    Reichel, M., Minchev, A., Larsen, N. L., 2014. Trim optimization-theory and
815        practice. International Journal on Marine Navigation and Safety of Sea

Transportation 8, 387–392.

Ronen, D., 2011. The effect of oil price on containership speed and fleet size. Journal of the Operational Research Society 62(1), 211–216.

Siroky, D. S., 2009. Navigating random forests and related advances in algorithmic modeling. Statistics Surveys 3, 147–163.

Shao, W., Zhou, P., Thong, S. K., 2012. Development of a novel forward dynamic programming method for weather routing. Journal of Marine Science and Technology 17(2), 239–251.

Sherbaz, S., Duan, W., 2014. Ship trim optimization: assessment of influence of trim on resistance of MOERI container ship. The Scientific World Journal 2014, 1–6.

Soner, O., Akyuz, E., Celik, M., 2018. Use of tree based methods in ship performance monitoring under operating conditions. Ocean Engineering 166, 302–310.

Song, D. P., Li, D., Drake, P., 2015. Multi-objective optimization for planning liner shipping service with uncertain port times. Transportation Research Part E: Logistics and Transportation Review 84, 1–22.

Takashima, K., Mezaoui, B., Shoji, R., 2009. On the fuel saving operation for coastal merchant ships using weather routing. In: the 8th International TransNav Symposium Proceedings, 431–436.

Townsin, R. L., Kwon, Y. J., 1993. Estimating the influence of weather on ship performance. Proceedings of International Symposium on TransNav 1993, 191-209.

Tsou, M. C., Cheng, H. C., 2013. An ant colony algorithm for efficient ship routing. Polish Maritime Research 20(3), 28-38.

Wang, K., Yan, X., Yuan, Y., Li, F., 2016. Real-time optimization of ship energy efficiency based on the prediction technology of working condition. Transportation Research Part D: Transport and Environment 46, 81–93.

Wang, S., Meng, Q., 2012. Sailing speed optimization for container ships in a liner shipping network. Transportation Research Part E: Logistics and Transportation Review 48(3), 701–714.

Wang, S., Meng, Q., Liu, Z., 2013. A note on "berth allocation considering fuel consumption and vessel emissions". Transportation Research Part E: Logistics and Transportation Review 49(1), 48–54.

Wang, S., 2016. Fundamental properties and pseudo-polynomial-time algorithm for network containership sailing speed optimization. European Journal of Operational Research 250(1), 46–55.

Wang, S., Ji, B., Zhao, J., Liu, W., Xu, T., 2018. Predicting ship fuel consumption based on LASSO regression. Transportation Research Part D: Transport and Environment 65, 817–824.

Wang, S., Wang, X., 2016. A polynomial-time algorithm for sailing speed optimization

854       with containership resource sharing. Transportation Research Part B:
855       Methodological 93, 394–405.

856   Wang, Y., Meng, Q., Du, Y., 2015. Liner container seasonal shipping revenue
857       management. Transportation Research Part B: Methodological 82, 141–161.

858   Xia, J., Li, K. X., Ma, H., Xu, Z., 2015. Joint planning of fleet deployment, speed
859       optimization, and cargo allocation for liner shipping. Transportation Science 49(4),
860       922–938.

861   Yang, L., Chen, G., Rytter, N. G. M., Zhao, J., Yang, D., 2019. A genetic algorithm-
862       based grey-box model for ship fuel consumption prediction towards sustainable
863       shipping. Annals of Operations Research, 1-27.

864   Yao, Z., Ng, S. H., Lee, L. H., 2012. A study on bunker fuel management for the
865       shipping liner services. Computers & Operations Research 39(5), 1160–1172.

866   Zhao, J., Yang, L., 2018. A bi-objective model for vessel emergency maintenance under
867       a condition-based maintenance strategy. Simulation 94(7), 609-624.

868
869
870
871
872
873
874
875
876
877
878
879
880

**Appendix A. Construction of decision tree regressor (Friedman et al., 2001; Harrington, 2012; Breiman, 2017)**

The input information for decision tree construction contains the training set and termination conditions. We denote the set of $J$ input features as $(x_1, x_2, ..., x_J)$. An input feature is denoted by $x_j$, and the value range of this $J-$ feature is $[x_j^{min}, x_j^{max}]$. A specific value of this feature is denoted by $s_j$, $s_j \in [x_j^{min}, x_j^{max}]$. In addition, we denote the training set containing $N$ data entries as $D = \{(x^1, y^1), (x^2, y^2), ..., (x^N, y^N)\}$. A data entry is denoted by $(x^e, y^e)$ with $e = 1, ..., N$, where $x_e = (x^{e1}, x^{e2}, ..., x^{ej}, ..., x^{eJ})$ is a dimensional vector containing $J$ features and $y^e$ is a one dimensional output value. The construction process of a regression decision tree based on CART algorithm requires finding the *best split* pair $(j^*, s_{j^*})$, $s_{j^*} \in [x_{j^*}^{min}, x_{j^*}^{max}]$ of the nodes when splitting. Denote termination condition (a) to (c) as $T_a$, $T_b$, and $T_c$. The main steps to construct a CART decision tree are presented as follows:

*Procedure 1: Construction of CART decision tree*

| | |
|---|---|
| *Input* | Training set $D$ and termination conditions $T_a$, $T_b$, and $T_c$. |
| *Output* | Regression tree $f^{DT}(x)$. |
| *Step 1* | Find the *best split* pair $(j^*, s_{j^*})$ of the current splitting node by solving the following formula: $$(j^*, s_j) \in \arg\min_{\substack{j \in (x_1, ..., x_J) \\ s_j \in [x_j^{min}, x_j^{max}]}} [ \sum_{e \in R_{m1}(j, s_j)} (y^e - \frac{1}{|R_{m1}(j, s_j)|} \sum_{e_1 \in R_{m1}(j, s_j)} y^{e1})^2 + \sum_{e \in R_{m2}(j, s_j)} (y^e - \frac{1}{|R_{m2}(j, s_j)|} \sum_{e_2 \in R_{m2}(j, s_j)} y^{e2})^2 ],$$ where $R_{m1}(j, s_j) = \{e = 1, ..., N \mid x^{ej} \le s_j\}$ and $R_{m2}(j, s_j) = \{e = 1, ..., N \mid x^{ej} > s_j\}$. |
| *Step 2* | Use the *best split* $(j^*, s_{j^*})$ to split the current node into two nodes that contain two sub datasets $R_{m1}(j^*, s_{j^*}) = \{e = 1, ..., N \mid x^{ej} \le s_{j^*}\}$ and $R_{m2}(j^*, s_{j^*}) = \{n = 1, ..., N \mid x^{ej} > s_{j^*}\}$ with output values as $c_1 = \frac{1}{|R_{m1}(j^*, s_{j^*})|} \sum_{e_1 \in R_{m1}(j^*, s_{j^*})} y^{e1}$ and $c_2 = \frac{1}{|R_{m2}(j^*, s_{j^*})|} \sum_{e_2 \in R_{m2}(j^*, s_{j^*})} y^{e2}$, respectively. |
| *Step 3* | Repeat step 1 and step 2 in a depth-first manner until coming to a node that reaches one of the preset termination conditions. Then, this node becomes a leaf node and a new node for splitting is found by backtracking. |
| *Step 4* | Repeat Step 3 until there is no more nodes that can be split. Finally, the total training set is separated into $M$ mutually exclusive sub-sets $R_1, R_2, ..., R_M$, and a sub-set is denoted by $R_m$. The decision tree model can be presented by |

$$f^{DT}(x) = \sum_{m=1}^{M} c_m I(x \in R_m), \text{ where } I(x \in R_m) = \begin{cases} 1, x \in R_m \\ 0, x \notin R_m \end{cases}.$$