

## Data science and advanced analytics for shipping energy systems

Coraddu, Andrea; Kalikatzarakis, Miltiadis; Walker, Jake; Ilardi, Davide; Oneto, Luca

**DOI**

[10.1016/B978-0-12-824471-5.00014-1](https://doi.org/10.1016/B978-0-12-824471-5.00014-1)

**Publication date**

2022

**Document Version**

Final published version

**Published in**

Sustainable Energy Systems on Ships

**Citation (APA)**

Coraddu, A., Kalikatzarakis, M., Walker, J., Ilardi, D., & Oneto, L. (2022). Data science and advanced analytics for shipping energy systems. In F. Baldi, A. Coraddu, & M. E. Mondejar (Eds.), *Sustainable Energy Systems on Ships: Novel Technologies for Low Carbon Shipping* (pp. 303-349). Elsevier. <https://doi.org/10.1016/B978-0-12-824471-5.00014-1>

**Important note**

To cite this publication, please use the final published version (if applicable). Please check the document version above.

**Copyright**

Other than for strictly personal use, it is not permitted to download, forward or distribute the text or part of it, without the consent of the author(s) and/or copyright holder(s), unless the work is under an open content license such as Creative Commons.

**Takedown policy**

Please contact us and provide details if you believe this document breaches copyrights. We will remove access to the work immediately and investigate your claim.

***Green Open Access added to TU Delft Institutional Repository***

***'You share, we take care!' - Taverne project***

**<https://www.openaccess.nl/en/you-share-we-take-care>**

Otherwise as indicated in the copyright section: the publisher is the copyright holder of this work and the author uses the Dutch legislation to make this work public.

## CHAPTER 7

# Data science and advanced analytics for shipping energy systems

Andrea Coraddu<sup>a</sup>, Miltiadis Kalikatzarakis<sup>b</sup>, Jake Walker<sup>a</sup>, Davide Ilardi<sup>c</sup>, and Luca Oneto<sup>c</sup>

<sup>a</sup>Faculty of Mechanical, Maritime, and Material Engineering, Delft University of Technology, Delft, the Netherlands

<sup>b</sup>Department of Naval Architecture, Ocean & Marine Engineering, University of Strathclyde, Glasgow, United Kingdom

<sup>c</sup>Department of Informatics, Bioengineering, Robotics and Systems Engineering, University of Genova, Genova, Italy

### 7.1. Data availability

The use of sensor technologies is rapidly expanding in the shipping industry, allowing for real-time monitoring and control of systems and processes. It is then possible to say that maritime data analytics is embracing the process of datification. Lloyd's Coffee houses and the publication of *Lloyd's list* in 1734 listing vessels and cargoes arriving in the port of London is one of the first records of data analytics in the maritime field. Nonetheless, the first big step in gathering and analyzing digital information is dated in the early 1990s. The first edition of [1], along with the establishment of *Shipping Intelligence Weekly* by Clarkson Research, powered a new generation of shipping analysts to develop tools that investigated changes in market cycles and vessel demands. Currently, the industry is hesitantly entering the third *Age of Maritime Data Analytics*, with the use of algorithms integrating several strands of data, from a component level to a fleet-wide level.

A modern seagoing vessel can generate a significant amount of data in a large variety of formats, which can provide an analyst with an holistic view of the vessel in terms of both internal and external awareness. In this context, internal awareness refers to all information regarding events occurring within the vessel (*endogenous information*), whereas external awareness (*exogenous information*) provides insight about the interaction of the vessel with its surrounding environment.

#### 7.1.1 Datification

Technological progresses made substantial steps forward in the last decades. Datification, namely, the process of transforming a phenomenon into data using sensors, is one of the fields that has most benefited from these technological evolution. Daily life of every individual is monitored by smartphone, smartwatches, and home automation devices. The industries are full of embedded devices to monitor the production processes real-time. Products produced by industry are natively equipped with sensors able to monitor

their own status. Moreover, an increasingly high number of sensors are installed on assets like bridges, wind turbines, highways and ships. All of these additions are motivated by the intent, at a different level, to predict the future whether to avoid adverse events or to profile users behavior.

The *datification* concept is not so novel. In 3800 B.C., the King of Babylonia recorded the first census in ancient Mesopotamia [2]. It was the first mean to measure the richness and powerfulness of the kingdom at that time and the very first datafication process ever made. Since then human beings never stop datifying the world all around. Nevertheless, the gap between the amount of past and current collected data is enormous. As Google CEO Eric Schmidt stated that [3]: “Every two days now we create as much information as we did from the dawn of civilization up until 2003”. What is new is the process of datifying aspects of the world and of our daily life that would have never been possible in the past. For example, social networks transform individuals into a live stream of heterogeneous different data sources. Industrial data sources, that were previously ignored or discarded (e.g., activity logs and machinery signals), are now becoming a crucial element to empower competitiveness. This datification process allows to transform the decision-making processes from a matter of best practice, intuition, and experience into a measurable science.

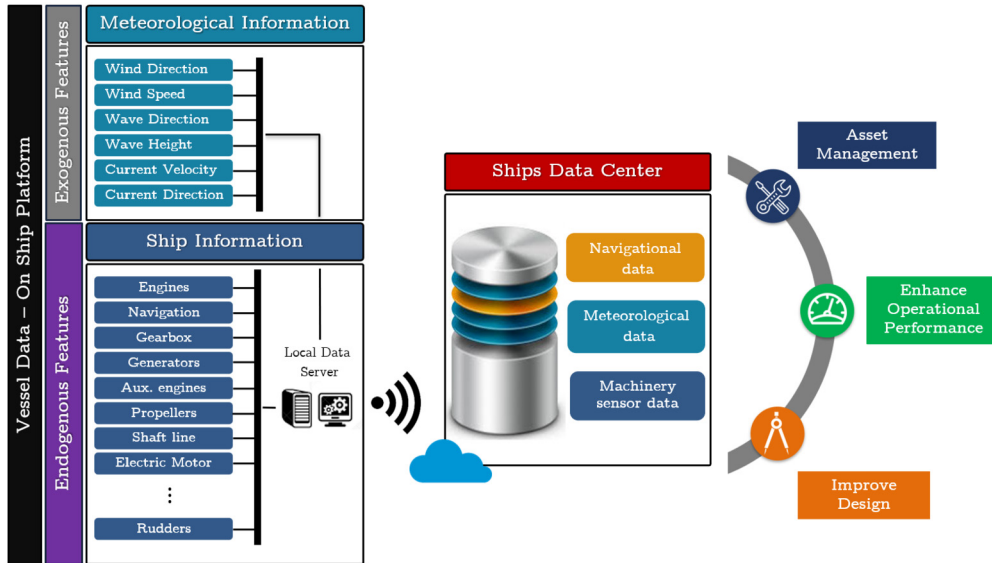
One crucial benefit of the process of datification is the ability to fuse together different data sources, namely the one directly linked to the phenomena under exam (the so-called endogenous sources) with the ones that can increase our knowledge base since they indirectly improve the endogenous sources (the so-called exogenous sources). A simple representation of a datification process of a vessel is reported in Fig. 7.1.

## 7.1.2 Endogenous data

Endogenous, in the shipbuilding industry, refers to the data generated by vessels during their life. A modern seagoing vessel can generate a significant amount of data in a large variety of formats. Automation and control systems, maintenance and condition monitoring, cargo monitoring systems, and equipment specifications all provide valuable information regarding the operational status of the vessel with respect to internal awareness. Let us briefly explore the various data sources and their scope.

### 7.1.2.1 Automation and control systems

As the market is driving shipowners to become more efficient with reduced manning requirements on-board, the use of sophisticated automation and control systems is becoming increasingly common and modern vessels are capable of safely operating for extended periods of time with unattended machinery spaces. To enable this, modern automation and control systems are fully integrated and capable of covering several aspects of the vessel's operations, and it is common for them to collect tens of thousands



**Figure 7.1** The datification process of a vessel.

of process measurements, and control system status information. In state-of-the-art systems, these aspects include the operation of the propulsion plant, power management of the electric power distribution systems, the operation of all auxiliary machinery, as well as navigation and administration of maintenance and purchasing of spares.

Monitoring and controlling the propulsion and the electric power generation systems is an essential task to maintain the efficiency and reliability of the vessel. There are numerous systems on-board that work together to provide electricity and propel the vessel, which require constant monitoring of different parameters such as fuel consumption, combustion and engine temperatures, over-load and over-speed limits, starting and stopping operations, generator voltage and frequency, load and torque of electric motors, heavy consumers' logic, and thruster monitoring.

The auxiliary machinery monitoring and control requires to supervise multiple systems and parameters. For instance, sea and fresh cooling water installations require pump and system pressure and temperature monitoring, potable and fresh water, bilge and sludge control require close monitoring of tank levels, pump pressures and valve status. Fuel oil system control requires monitoring and processing of tank levels, temperature, viscosity and flow, purifier and heater status. Understandably, similar complexity and amount of information arises from Heating, Ventilation and Air Conditioning (HVAC) systems, ballast water treatment, and exhaust gas temperature treatment systems.

The process of monitoring of all these systems generates a significant amount of data that can be extracted and used by analysts for various purposes, such as identifying instrument failures, finding and quantifying mechanical issues, measuring the effective-

ness of the control systems on-board, and identifying control strategies that are no longer working properly.

### **7.1.2.2 Maintenance and condition monitoring**

The reliability level required by modern vessels can only be ensured with timely maintenance, which involves periodic checks, repairs, and equipment replacements. All these events are routinely recorded, as required by the International Safety Management (ISM) code, in various forms that have evolved over the years.

In modern vessels, these events are cataloged digitally in maintenance management software, which provides an interface that enables seafarers and engineers to:

- schedule and document all planned and unplanned maintenance events;
- define and schedule time- and condition-based maintenance tasks;
- provide details regarding the criticalities of each event;
- automatically keep track of available spare parts;
- generate and keep track of life-cycle records of each equipment.

In addition to the data presented above, shipping companies routinely use additional sources when developing maintenance procedures for a particular vessel, which can further enrich the information available to engineers or analysts in terms of:

- maintenance guidelines given by the manufacturer;
- equipment history that includes defects, damages and remedial actions taken;
- equipment criticalities;
- age of the vessel;
- third party inspections;
- planned maintenance intervals;
- International Safety Management (ISM) guidelines.

Moreover, it is becoming increasingly common for shipping companies to use condition monitoring software (e.g., [4]) for the monitoring of the most critical components. Such software are constantly monitoring key parameters of the machinery using various methods (e.g., vibration and temperature measurements, ultrasonic signals, thermography, current analysis) to identify subtle changes that are indicative of developing faults. This information, when available, further supplements the rest of the data sources related to maintenance events. When properly analyzed, all this information does not only provide an overview of the life-cycle costs of the vessel, it can also help predict equipment failures in advance and allow shipowners the freedom to address developing defects in a way that best suits their operational goals.

### **7.1.2.3 Cargo monitoring**

Cargo monitoring and control systems are prevalent in oil and chemical tankers and Liquefied Natural Gas carriers and containerships. They differ based on vessel type and cargo and record different parameters, however, their scope always remains the same: to

ensure the safety of the crew, vessel and cargo, and facilitate efficient loading / unloading and storage.

Depending on the complexity of the process and the type of cargo, vessels may have a dedicated cargo control room to monitor important Key Performance Indicators (KPIs) of the state of the cargo, as well as all the systems that are involved in the loading / unloading processes and cargo storage. In oil and chemical tankers and Liquefied Natural Gas carriers that carry sensitive and dangerous cargo, these systems are logging and processing of cargo temperatures, pressures and flows of all cargo and ballast pumps, tank levels, heeling angle estimation and control, as well as constant monitoring of trim and list angle and draft. Similar parameters are monitored in containerships, including the position and status of gantries, hatch covers, gearboxes and special cargo (dangerous goods and reefer cargoes) are monitored to ensure the stability of the vessel and potential loss of cargo.

Modern systems also have the ability to benchmark every port call to identify efficiency improvements and reduce the risk of a delay, and even record berth performance so that analysts can compare the reported performance from the terminal, against the observed performance from the vessel.

#### **7.1.2.4 Equipment specifications**

Modern vessels are large and complex platforms that must be self-sustaining in their environment for extended periods with a high degree of reliability. There is a wide variety of components on-board that work in synergy to realize all the key functions of a seagoing vessel, such as machinery required for propulsion, steering, anchoring and ship securing, cargo handling, air conditioning and ventilation, and power generation and distribution. All these components are accompanied by various technical specifications and safety sheets, application guides, and user manuals that detail key characteristics of the machinery, underline functional information and KPIs, the intended scope of use, and maintenance specifications and instructions. This documentation provides a finer understanding of the equipment capabilities and makes the equipment's design, metrics and capacity clearly understood. Furthermore, it gives a detailed overview of the service conditions and processes that should be followed for proper equipment maintenance.

This information, as intended, helps engineers in establishing routine equipment conditions, usage frequency and the environmental conditions in which the equipment may or may not be used. However, it is very beneficial to analysts who can generate accurate evaluations of the equipment functioning under ideal operating conditions without any deviation. When these analyses are supplemented and compared with real-time information about the equipment of a sea going vessel, allow us to evaluate whether the equipment is appropriate to meet current and future objectives, derive an actual appraisal of the equipment's current efficiency and level of activity, and alarm us against costly mistakes that can disrupt operations, or cause material and human loss.

**Table 7.1** Example vessel data features.

<b>Vessel feature</b>	<b>Sensor type</b>
Position	Global Positioning System (GPS) Receiver
Speed Over Ground (SOG)	GPS Receiver
Speed Through Water (STW)	Doppler Log
Sink	Hydrostatic Pressure Sensors
Vessel motions	Motion Reference Units

### 7.1.2.5 Vessel environment interaction

Some features are usually utilized to describe and understand the interaction between the vessel and its surroundings and do not fit strictly into the endogenous / exogenous categorization. For the purposes of this chapter, they are considered as a subset of endogenous features. Table 7.1 describes the data features that measure the interaction between the vessel and its environment.

The Global Positioning System (GPS) is one of the most important devices for measuring vessel's position. GPS work by relating the position of the receiver to a constellation of GPS satellites. The receiver uses trilateration to determine the position of the vessel given the distance from the receiver to each of the 3 most proximate satellites in the network. The GPS receiver is responsible for tracking the vessel position throughout a voyage. In addition to serving a navigational role, the change in position is used to determine the GPS speed, which is also the Speed Over Ground (SOG). The position and speed of the vessel are logged into the vessel data center also accounting for the timestamp.

Speed Through Water (STW) is the headway speed due to the force produced by the vessel. In deep water scenarios, the STW of a vessel is determined by the Doppler Log. The Doppler log measures the STW by emitting a signal at a known frequency from the bottom of the vessel, which reflects off the sea bed and detected by the device at the new frequency. The Doppler log then uses the Doppler Shift equation to infer the STW by measuring the relative shift in frequency between the source signal and the apparent one. Configurations of the Doppler log with only one transducer are often sensitive to the transient motion of the vessel such as slight pitch and roll. To overcome these sensitivities, the Janus configuration is used to average the results of four signals for a more robust prediction of the speed. In shallow water, most Doppler logs will only be capable to measure the SOG. However, recent advancement in the JLN log from Japanese Radio Company (JRC) [5] is able to measure the STW even in shallow depths (2 meters). It should be noted that the speed of sound in water varies between 1,450 and 1,570 meters ( $\sim 9\%$ ) depending on salinity and temperature, which can induce errors in logs which do not consider these parameters when calculating the speed. However, the International Maritime Organization (IMO) performance standards for the Doppler log mandate the performance for depths greater than 3 meters from the keel. The displayed



measurements for the speed and distance through water must be within the greater of either, a 2% or 0.2 knots tolerance for digital displays (or equivalently 2.5% or 0.25 knots for analogue displays) even in conditions where the vessel roll is up to 10° & 5° pitch.

To measure the vessel motions and displacement requires a combination of sensors working in tandem to capture the dynamic behavior. The draft is usually measured from a set of hydrostatic pressure sensors deployed at the bow and stern. With larger ships, it is beneficial to include additional sensors at the port and starboard mid-ship. From the hydrostatic pressure readings, it is then possible to infer the draft along the length of the vessel and negate temporary wave effects. It is possible to compute the trim from the difference between the fore and aft drafts and the length of the ship. This approach yields a good approximation of the trimming angle when stationary, but is not as robust in sailing conditions. On the other hand, specialized sensors to detect the exact position of the fore and aft of the ship are also deployed to measure the Dynamic Trim in transient conditions. In this case, two GPS sensors are deployed at the fore and aft masts which continuously monitor the dynamic trim regardless of the sailing condition. Finally, the vessel motion response is usually recorded with an Inertial Measuring System (IMS). Similar to the other dynamic responses, two IMSs are deployed at the fore and aft of the vessel to ensure a robust measurement of the vessel motion. The IMS contains accelerometers and gyroscopes to measure the motion and rotation in 6 degrees of freedom.

#### **7.1.2.6 Navigational data**

All this equipment is a valuable source of information regarding external awareness, as there is a significant amount of data generated from the sensors and relevant measuring equipment found on the bridge, such as radars, rate of turn indicators, heeling angle recorders, the Electronic Chart Display and Information System (ECDIS), auto-pilot, and the Automatic Identification System (AIS), which transmits the vessel's unique identification number, position, speed and course as required by the IMO for all commercial vessels over 300 Registered Tonnage. Moreover, additional instrumentation and data sources can be found on the bridge of special-purpose vessels, which can also include wave radars, oil spill detectors, and high-accuracy inertial navigation sensors. These data sources, when combined with all data providing us with internal awareness, can give a holistic view of a seagoing vessel and its surrounding environment.

#### **7.1.3 Exogenous data**

Environmental and navigational data provide us with information regarding the surrounding environment in which the vessel is operating. Historically, weather forecasts and navigational data were captured through descriptive notes, or later transmitted via

**Table 7.2** Exogenous data.

<b>Climate features</b>	<b>Sensor type</b>
Wind speed and direction	Anemometers (Mechanical or Ultrasonic)
Air Temperature	Thermocouples
Relative humidity	Hygrometer
Pressure	Barometers
<b>Metocean features</b>	<b>Sensor type</b>
Current speed and direction	Inferred with Doppler Log
Seawater temperature	Thermocouples
Sea depth	Echo Sounders
Significant wave height	Satellite Radar Altimeters

radio, were the primary source of data for seafarers to draw inferences regarding the likelihood of the success of a voyage. Nowadays these data can be transmitted both using satellite communication (in open sea) and VSAT, L-band, and 3G/4G/LTE networks in coastal navigation. Table 7.2 describes the exogenous data that can be collected.

### 7.1.3.1 Climate features

At the present time, weather satellites are one of the primary source of climate data utilized in the maritime industry. This technology is widely considered the standard in weather feature data collection due to the variety and accuracy of the data gathering methods. The process of collecting these data is explored in [6]. Weather data from these satellites is often gathered and distributed for free by many state-funded organizations, such as the U.S. Naval Research Laboratory (NRL) [7], National Weather Service (NWS) [8], and the MET Office [9], although industrial services exist to collate the data and provide access to maritime clients at varying degrees of detail [10–12]. Radio occultation technology is also at the forefront of climate data capturing. A competitive proponent of this method proposed in [11] aims to improve the accuracy of weather forecasting with less reliance on calibration. They provide an exhaustive option in regards to maritime forecasting as their platform incorporates a high frequency network of Global Positioning System (GPS) / Global Navigation Satellite System (GNSS) satellites and emphasize learning from previous occultation data to improve performance.

State-of-the-art research into the datafication of weather features is driving the accuracy and availability of these datasets. Here we discuss technologies applied to capture the most important climate and metocean features and the principal implementation scenarios. Robust environmental data loggers used to capture this information include considerations towards both climate and metocean attributes. Now, like many other maritime processes, advancements in deployable, economic sensing technology have led to a transformation in how we apply this data in the context of marine energy systems.

From a mariner's perspective the weather data, specifically regarding the wind, is essential for understanding the stability of a vessel. Especially in the case of container vessels, VLCCs, LGN, and Liquefied Petroleum Gas (LPG) vessels due to the wide side ("sail") area. The force from the wind also influences the surface tension force in the ocean that causes swell and variability in the sea state. In combination with data describing the ocean current, this is a powerful inference on the length of a voyage, the vessel's fuel consumption, and the safety during operation.

Although the wind has arguably the most direct effect on the sailing conditions, it exists only due to differences in atmospheric pressure as air flows from high to low pressure areas. In fact, historically, recording the pressure was one of the most important inferences mariners had about the short term weather forecast. Other climate factors that influence the sailing conditions include the humidity which describes the ratio of water vapor in the air, and the air temperature. Ultimately, the weather is the direct product of the transfer of energy in the environment. When the different sources of climate data are combined it is possible to develop a clear understanding about how the weather will influence the success of a voyage.

### **7.1.3.2 *Metocean features***

In order to measure waves, ocean wave monitoring buoys offer an alternative method of datafication for weather features and are employed in applications where satellites are not suitable or available. Floating weather stations have been in operation from the early 20th century where ships designated for the datafication of weather features were anchored to form a network of weather stations. Today, ocean buoys are deployed with sensors to capture both climate and metocean features, in comparison to satellite technology, the data collection from a buoy is able to effectively target an area of interest with a cardinality restricted only by the frequency of the sensors rather than the period of orbit. However, the available region the buoy can cover is minute with respect to the satellite (typically in the region of two orders of magnitude [13]) and is usually reserved for instances where the primary method is unsuitable. State-of-the-art reviews of weather systems over the last decade [14,15], have centered around improving the performance of X-band radar (8.0 to 12.0 GHz) for capturing metocean conditions by analyzing the sea surface spectra. The ability of the X-Band system to assume the calculation of these features is of extreme value in the marine industry as it comes at almost no additional cost to the vessel, since radar systems are almost universally employed in vessels for a suit of other applications, mainly: for safety purposes so as to avoid collisions and navigation purposes. The most prominent application of this technology, wave motion detection, has been applied through a range of frequencies to accurately describe the changing sea state and vessel response. Research has shown the expansion of this technology from the defense industry to the commercial sector has proven its use as a marine remote sensing tool within a range of several kilometers when attached to a moving vessel.

A comparison between short and medium wave inferences in [16] has shown that the significant wave height estimated by the medium wave pulse width is in agreement with the short wave after a re-calibration process, further research in [17] then demonstrates the implementation of this process without disrupting the essential safety mechanisms performed by the X-Band system.

For what concerns the currents, the information describing them, is often implicit within other data sources in shipping applications. For example, by comparing the difference between the SOG and STW the effects of the currents can be obtained. For problems that are specific to a location, the Acoustic Doppler current Profilers (ADPs) work on the same principle as described in Section 7.1.2.5. However, these devices (which are discussed further in [18]) are moored as buoys and can often measure both the currents and waves in a single device.

With consideration towards mapping ocean depth, shore-based radar stations have long been able to map the depth of the sea bed in shallow water, without the need to travel directly to the location. With sufficient georeferencing data it was shown in [19] that X-Band radar technology, mounted on a vessel, could favorably determine the bathymetry up to depths of 50m kilometers away from the points of interest. This presents a potential safety mechanism when working in areas of uncharted water, since the radar can effectively determine ocean depth in areas that vessels may not be able to enter.

#### 7.1.4 Discussion

Obtaining a holistic view of a seagoing vessel requires aggregating tabular, unstructured, and geospatial data, with sampling rates that span several orders of magnitude, from multiple and heterogeneous data sources, generating datasets of different structure. Even though datafication is a reality in several sectors (e.g., finance, media, telecommunications, and healthcare), its adoption in the maritime industry has been slower, and the full benefits have not fully materialized yet. Slowly but surely, the shipping industry will evolve, from using a decision-tree driven, to a data-driven approach.

### 7.2. Data science and advanced analytics technologies

Data Science and Advanced Analytics are the fields of research that study how to exploit data to derive new, meaningful, and actionable information [20–35,76]. Since in this chapter we are focusing on Shipping Energy Systems, in this section we will describe the subset of methods, technologies, and tools that can exploit the data listed in Section 7.1 to answer the key questions of this domain. In particular, Data Science and Advanced Analytics allow to answer four fundamental questions.

The first question is “*What happened or is happening to the system?*”. This question, even if naive at first sight, requires at least a datafication process, namely, to have properly

sensorized the systems, collected the data produced by the sensors in a data hub, and designed the required dashboards and KPI to be able to visualize the performance of the systems and then their status. This process is not a trivial one since, in the shipping context, there is still a large space for human intervention: the status of the engine is manually checked by experienced operators, some sensors cannot be deployed in legacy systems, the data are not centrally collected, the data suffer of low quality or reliability, the right KPI are not easy to be synthesized, and the dashboards are not informative enough to be fully reliable. Even if in the last years a large effort has been spent in order to fill these gaps, there are still some major issues to be addressed: exploit the free form or manual reports of the operators, build virtual sensors able to estimate sensors measurements that cannot be deployed in legacy systems, build a network infrastructure able to centrally collect all the data produced, design a limited set of meaningful and informative KPI, build dashboards able to help the user in easily and rapidly understanding the system status. *Descriptive Analytics* is the field of Data Science and Advanced Analytics that allows to answer these questions and leverage the technologies coming from the world of data collection, storage, cleaning, cleansing, exploration, and visualization, namely Data Wrangling.

The second question is “*Why is something happening or has happened to the system?*”. This second question is clearly more complex than the first one and requires to extract information that is not readily available even if all the data regarding the system are. Especially in the Shipping Energy Systems domain, this question is commonly answered by experienced operators who gained, during many years of direct experience on the field, the required knowledge to be able to understand, based on data or manual/visual inspections, why a particular problem arose in the system. In fact, most of the problems of a Shipping Energy System are recurrent and very few of them are novel problems never happened in the past. For this reason, the experience is essential since it allows to rapidly detect the reasons of these repetitive problems. Data allow to answer this question in a more scalable and reliable way since one single operator cannot see, in its entire life, a hundred of different systems while computers can store the information of possibly all the systems for their entire working life. Therefore, it is possible to exploit all these data to find correlations and patterns which allow to understand in which situations, and then why some problems arise. Note that correlation and causality are similar concepts but not the same ones (causality implies correlation but not vice versa since spurious correlations are widespread in nature<sup>1</sup>). Nevertheless, listing all the correlated factors is surely an effective way to restrict the possible causes. *Diagnostic Analytics* is the field of Data Science and Advance Analytics that allow to answer these questions and leverage the technologies coming from the world of Data Mining.

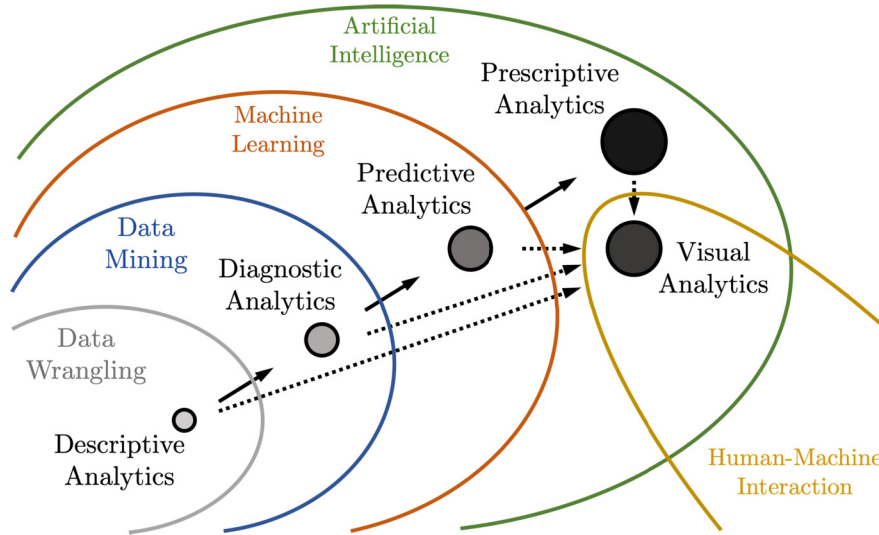
The third question is “*What will happen to the system?*”. Answering this third question is even more challenging with respect to answering the previous ones but surely provides

<sup>1</sup> <https://www.tylervigen.com/spurious-correlations>

a first big insight toward deeply understanding the behavior of the system, and it may allow to act before something can happen. The answer to questions about the future is always associated with a risk, namely, the probability of mistakes associated with a prediction, which is quite complex to estimate. In fact, Descriptive and Diagnostic Analytics, which pose questions on present or past system behaviors, require answers whose risks can be easily estimated using past data. Instead, when we pose questions about the future, we need to make some hypotheses and be completely aware of their meaning to make aware decisions that fully contemplated these risks. As an example, let us suppose that we want to migrate the maintenance processes of a component of a ship energy system from a corrective or working-hours-based maintenance to a condition based or predictive maintenance. In this scenario we will have to build models able to estimate and predict the current and future decay status of this component and, for this purpose, we need to have some data regarding its past behavior. Once these models have been built, we have to keep in mind some important considerations. For example, the first one is about the representatives of the exploited data (namely, the ship will be still used in the future for the same purposes? Or will it travel the same routes?). The second one is about the possible causes of malfunction. If in the past the component has been maintained with a higher frequency to avoid any possible issue we have to take into account that the model has never seen data regarding very high decay state or malfunctions so it will probably never be able to extrapolate too much. Consequently, answering questions about future behavior needs to take into account many complex relations between systems, data, models, decisions, and actions. *Predictive Analytics* is the field of Data Science and Advanced Analytics that allows to answer these questions and leverage the technologies coming from the world of Machine Learning (ML).

The last question is “*What is it needed to do to make the system behave in a certain way?*”. This last question is surely the most challenging one and requires a lot of technologies since it points out to the design of a fully automated supervision system able to guide the ship energy system toward a particular behavior. For example: why do I have to warm-up an engine to reduce the maintenance? Why do I have to modulate the trim to reduce the energy consumption? How often do I have to clean the hull to optimize the trade-off between maintenance costs and fuel savings? In order to address these questions, we do not have to simply predict the behavior of a system, but we also have to model, for example, the system constraints, the preferences of the operator which exploits the system, the maintenance constraints, the maintainers preferences. For this purpose, we need all the information listed in Section 7.1 and we need to be able to both build data-driven models, describe the knowledge bases, describe the constraints, and describe the processes. *Predictive Analytics* is the field of Data Science and Advanced Analytics that allows to answer this question and leverage on the technologies coming from the world of Artificial Intelligence.

Aside from these four questions, there is a field of Data Science and Advances Analytics called *Visual Analytics*, that exploits all the previously mentioned analytics



**Figure 7.2** Analytics approaches for shipping energy systems.

(Descriptive, Diagnostic, Predictive, and Prescriptive) tools and results to give visual insight to the human operators and let them make the best out of it. In fact, it is quite rare, also in Prescriptive Analytics, to let an intelligent system take a final decision on a complex subject in a fully autonomous way (e.g., decide to stop a vessel and make an important maintenance intervention to a ship energy system). The final decision is often taken by a human operator and Visual Analytics is the science which studies how to present the information coming from different analytics in the most possible aware manner. For this reason, a cornerstone of the Visual Analytics is the concept related to the Human–Machine Interaction, which allows to design ways to make the intelligence of the machine be fully exploited by the intelligence of the humans. Machines are often not very smart but very fast in doing simple things (the so called narrow intelligence), while humans are often slow but very clever in developing complex connections and solutions (the so called general intelligence). Visual Analytics tries to empower the human mind with the tools of Descriptive, Diagnostic, Predictive, and Prescriptive Analytics making them easily and fully exploitable by human for taking faster and more aware decisions and actions.

A graphical representation of the field of Data Science and Advanced Analytics and related technologies is depicted in Fig. 7.2.

Note that this segmentation of the Data Science and Advanced Analytics research field is not always so neat and it is not the only possible one. In fact, in some cases even, for example, for Descriptive Analytics, it is required to use ML technologies (e.g., to create a virtual sensor for a system that cannot be retrofitted with a particular

sensing technology as we will see in this chapter). Note also that, as we will see in this section, the distinction between, for example, Data Mining and ML or ML and Artificial Intelligence and Deep Learning, is again not always so clear and neat but this chapter presents the point of view of the authors acquired in many years of research in this field.

Once Data Science and Advanced Analytics aim has been understood, we can describe the technologies on which they leverage. We will start with the Data Wrangling, we will continue with Data Mining, then ML and finally with Artificial Intelligence.

Data Wrangling [36–40] is a broad term that contemplates many technologies, but in general, it allows to transform raw data into a manageable format that can be exploited to perform higher level analysis. We will use this concept with its broader meaning including data collection, storage, cleaning, cleansing, exploration, and visualization. Recent technologies for big data collection and fruition systems allow to collect, store, and access huge amounts of data from different and heterogeneous sources. Software frameworks for centralized and distributed data storage and processing (like Hadoop, Spark, Hive, MongoDB, etc.) and their ecosystems allow to easily access data in different formats from different sources, to create and to process big data sources with reference to the particular phenomenon under exam. Usually data can appear in nature in many different forms (see Section 7.1), from the classical tabular data, to images, free form reports, natural language, graphs, and they can represent different information like sensor logs, maintenance reports, configuration data, concepts, and relations. For this reason, it is necessary, and extremely useful, to collect it in a centralized hub, curate it, keep it always updated since data represent an intangible, but fundamental asset for any modern industry, especially the shipping one, that will become even more crucial in the future. In fact, progresses in shipping energy systems technologies are experiencing a sort of advancements plateau (like in many other industries), while the space for improvements in the adoption of data-driven technologies and solutions is still very large and in many fields of shipping it represents the one in which companies build their competitiveness. Data Wrangling allows to understand what happened and is happening to a ship energy system since it represents the source of information that needs to be synthesized in KPI to be displayed to human operators and to stakeholders to take decisions based on data and not just the experience.

Data Mining [41–45], contrarily to Data Wrangling, is a technology which allows to extract additional information from data and not simply collect, store, clean, cleans, explore, and visualize it. Data Mining main objective, as described before, is to understand why something has happened to a system based on the current and historical data. In other words, Data Mining main focus is to discover repetitive patterns related to a particular fact or behavior. The first step toward addressing this goal is to find a correct notion of similarity and distance between data and patterns. This may appear a simple procedure for structured data, but for text, graphs, concepts, and relations this

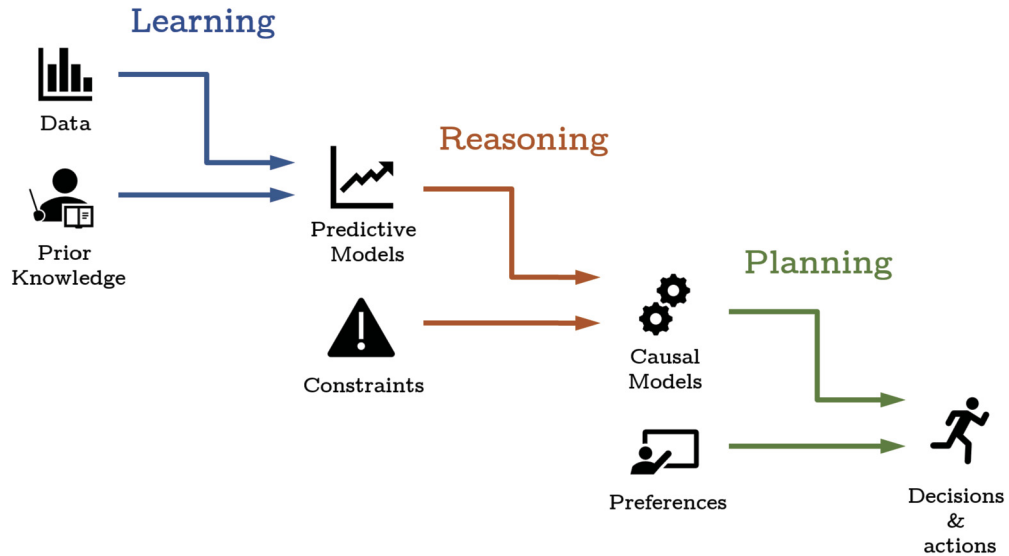


is not trivial at all. One has to define a metric of similarity or distance which should simultaneously well represent the phenomena, be rigorous enough to be mathematically described, and be simple enough to be efficiently computed. The second step in data mining is to apply algorithms able to find recurrent patterns or associations between them. For this purpose, we have many tools aiming at solving different problems: from Frequent Pattern Mining Model to Association Rules and Interesting Pattern Analysis with their different algorithms. There are also more advanced tools and concepts like Pattern Summarization and Pattern Querying, which are fundamental to make a synthesis when the number of patterns or amount of data increases. Finally, there are also ways to group data and patterns to find similarities at higher hierarchy. Clustering Methods and Algorithms focus on this problem. Many applications in shipping energy systems require the partitioning of data points into intuitively similar groups. The partitioning of a large number of data points into a smaller number of groups helps greatly in summarizing the data and understanding it for a variety of data mining applications.

ML [46–53] is a subset of Artificial Intelligence defined as the study of computer algorithms that improve automatically through experience. ML algorithms build a mathematical model based on sample data, known as training data, to make predictions or decisions without being explicitly programmed to do so. ML algorithms are able to solve two main types of problems, namely, classification and regression, and these problems can be faced in supervised, semi-supervised, and unsupervised fashion. The distinction between classification and regression lays in the output that the model has to predict. For classification we need to estimate the class belonging (so there is no concept of distance between the classes) of a particular observation e.g., when one has to detect if a maintenance is needed or not based on the data coming from the automation system, while in regression we have to estimate an output of a system where the possible output is a sorted set of possibilities (so there is a concept of distance between the outputs), e.g., when one has to detect the fouling status exploiting data coming from the automation system. Note that many problems that seem to not fit in this framework, like the time series forecast, e.g., when one has to detect the fuel consumption of the propulsion system based on the history of fuel consumption, can be easily plugged in a regression framework with an autoregression approach or appropriate decomposition. For what concerns, instead, the distinction between supervised, semi-supervised, and unsupervised cases, in the supervised case we have a possibly large amount of historical labeled data as training set, namely both the observations of the system and associated outputs are available. This case is the simplest one, but in many applications (e.g., fault detection) the number of faults is quite limited, or not always available, so the amount of labeled samples is limited. In these cases, called semi-supervised, we have a large amount of unlabeled samples while few data are labeled. Finally, the unsupervised case, which is the most complex one, consists in all those cases where no desired output is available for the observations (e.g., detect the operational profiles of a ship). Based on

the data available and the problem that one has to solve, it is possible to derive specific sub-problems, e.g., clustering is an unsupervised classification and novelty detection is a semi-supervised classification problem. For each specific sub-problem there are then families of algorithms (e.g., Ensemble Methods (EM), Kernel Methods (KM), Artificial Neural Networks (ANN), Gaussian Processes (GP), Random Projections, and Rule Based Methods) and then specific algorithms (e.g., Support Vector Machines (SVM) in KM or Random Forests (RF) in EM) that can be exploited. These algorithms can be easily divided in two groups: shallow and deep models. Shallow models require a so called feature engineering phase that maps the raw data (e.g., the data coming from the automation system) into a feature vector, called representation vector, using the knowledge about the physical problem (e.g., using frequency analysis) and then, by having a training set of features vector and possibly the associated output, the shallow model learns the relation between input and output. In many cases this feature engineering phase is not trivial (e.g., image of graph analysis) but a large training set of raw data and possibly the associated output is available. In this case, deep models are able to automatically learn from the data itself the feature vector, also called data representation or embedding, and the final model. Deep models are quite powerful but often require a huge amount of data to be trained and not all the applications have these amount of data available. Another issue that we have to face in ML is how to tune and assess the performance of the algorithms [54]. In fact, ML algorithms often require to choose between different approaches and each approach is characterized by hyperparameters to tune. These choices deeply affect the performance of the final model and for this reason they must be tuned and assessed carefully. Resampling techniques like cross validation and non-parametric bootstrap are often used by practitioners because they work well in many situations. Other alternatives exist, which are based on the Statistical Learning Theory, which give more insight into the learning process. Examples of methods in this last category are: the seminal work of the Vapnik-Chervonenkis Dimension, its improvement with the Rademacher Complexity, the theory of compression, the Algorithmic Stability breakthrough, the PAC-Bayes theory, and more recently the Differential Privacy theory.

Artificial Intelligence [55–59], in the context of this chapter, is the effort of fully, or at least partially, automatize the process of taking decisions. For this purpose, we need different tools able to build all the intelligence needed to start from data and design decisions and actions to undertake (see Fig. 7.3). A first building block is ML able to exploit data and prior knowledge about a particular domain to build predictive models. These ML models scale well with the amount of data available but they are not as effective if exploited for deduction purposes. For this reason we need a second building block which exploits the Model Based Reasoning technology. Model Based Reasoning (based, for example, on classical optimization or on Answer Set Programming, or in AI-guided optimizers) allows to model in an effective way complex systems and its



**Figure 7.3** Artificial intelligence for shipping energy systems.

constraints based on the physical knowledge about them and to deduce meaningful information by solving complex (optimization) problems creating causal models. The Model Based Reasoning limitation is that it may not scale well with the size of problem. The joint use of ML and Model Based Reasoning allows to fill their gaps and empower them with their strengths. Note that multiple solutions may result in the achievement of the desired outcome, we need also to model the preferences of the domain and use a preferred planner (e.g., using Planning Domain Definition Language) to finally take the final decision and actions.

Note that, because of the particular domain of the shipping energy systems, not always a fully automated system can be exploited or people will accept its use. In fact, some decisions require a final, or in between, supervision of a human expert or stakeholder. For this purpose, in many cases, just part of the building block depicted in Fig. 7.3 is actually exploited and then the information is provided, in some form, to a human expert. The way in which this information is presented requires a profound knowledge of technological and psychological aspects related to the Human Computer Interaction [60–65]. In fact, providing the information alone does not guarantee that the user will be able to exploit it. The way in which the information is presented is as-much or more important than the information itself. Human Computer Interaction enhances the capabilities of the information receivers by developing data- and context-driven interfaces. The interaction between decision-makers and the context-driven human machine interface allows to overcome the effects of missing or inaccurate measures and data and of model uncertainties. At the same time, decisions can be im-

proved by means of better situational awareness, helping the decision-maker to evaluate trade-offs in finding the best compromises toward the final decision. Particular attention needs to be given in avoiding a visual information overload to the decision-makers to provide the minimum number of stimuli with the maximum informative value. The stakeholders can then be able to accept the prescribed decisions or to further speculate by understanding the process that led to the prescribed suggestions. Furthermore, decision makers need to be able to check for alternative solutions by parameter steering or varying the context-specific constraints, preferences, and KPIs before taking the final decision. This continuous interaction between the system and the decision-makers should be also recorded, allowing to close the human-in-the-loop cycle and back-propagate the human corrections to system suggestions.

### **7.3. State-of-the-art in data science and advanced analytics**

In the following, the authors provide an overview of the state-of-the-art contribution in the context of data science and advanced analytics for shipping energy systems. We selected the most relevant (according to the venue) and referenced works in the last 6 years. These results are summarized in Table 7.3.

Lu et al. [66] address the problem of the energy costs of shipping by the development of a route optimization framework to determine the Energy Efficiency of Operation (EEO) of a vessel considering the variable resistance of a ship's hull. The proposed framework combines the ship characteristics determined from historical data and a modified Kwon's method for semi-empirical performance estimates, and corrections for the variable resistance of the vessel hull. The resulting model was able to predict the ship's resistance from Vessel – Environment interaction data (the encounter angle between the ship and the wind, and the draft) and exogenous data (mainly, the sea states) for the two case studies considered. The results demonstrated that the inclusion of these data sources led to increased accuracy in the prediction of hull resistance compared to traditional semi-empirical models which was then used to determine the EEO for the vessels. The Modified Kwon's method was also applied in a route optimization framework where the authors obtained routes for both the lowest fuel consumption and the lowest Beauford number & fuel consumption together. The authors admit the variability in both the accuracy of exogenous data forecasts and the vessel's ability to follow the prescribed route may lead to considerable uncertainties in the application.

Wang et al. [67] proposed a dynamic control system to identify effective energy saving and CO<sub>2</sub> reducing measures for a cruise ship. The authors proposed reducing the energy consumption by tuning a control system to predict the total resistance of the vessel in terms of parameters describing the Vessel – Environment interaction and suggesting an optimal trajectory. The control system was able to predict the Energy Efficiency Operational Indicator (EEOI) by solving a non-linear dynamic Time-Series

**Table 7.3** Summary of the state-of-the-art contributions in the context of data science and advanced analytics for shipping energy systems. We selected the most relevant (according to the venue) and referenced works in the last 6 years.

Ref.	Scope	Input data	Models output	Methods	Results
[66]	Route Optimization	Vessel – Environment Interaction & Exogenous Conditions	Route Fuel Consumption [t]	Semi-Empirical model	Accuracy not reported Reduction in Fuel Consumption of 7–11% for 4 routes
[67]	Control System Optimization	Vessel – Environment interaction & Exogenous Conditions (Forecasted)	Optimal set of control parameters to minimize EEIO	Semi-Empirical models PSO	Accuracy not reported Improvement in EEIO of ~2% Reduction in fuel consumption and CO <sub>2</sub> of ~28%
[68]	Fuel Consumption Reduction	Vessel – Environment interaction & Exogenous Conditions	Predicted Navigation Environment	Wavelet ANN	Accuracy not reported Reduction in fuel consumption per unit distance of 19.04%
[69]	Fuel Consumption Estimation	Vessel – Environment Interaction	Fuel Consumption [Mt]	ANN	MAPE: 7.4–10.8%
[70]	Fuel Consumption Estimation	Vessel – Environment Interaction & Exogenous Conditions	Fuel Consumption [L/h]	Bag, RF, Boo	RMSE: 45.2 (Bag) RMSE: 43.5 (RF) RMSE: 41.3 (Boo)
[71]	Operational Pattern Classification	Automation and Control Systems, Equipment Specifications, Metocean Features	Operation Type	SVM	Accuracy: 98% (port) Accuracy: 94% (sailing)
[72]	Shaft Power Estimation	Automation and Control Systems, Vessel – Environment Interaction, Climate Features	Shaft Power [kW]	ANN	MAPE: 7.8%
[73]	Wear fault diagnosis of Marine Diesel Engines	Automation and Control Systems, Maintenance and Condition Monitoring	Wear Fault Mode	BRB, ER, ANN	Accuracy: 86.7% (BRB) Accuracy: 83.3% (ER) Accuracy: 93.3% (ANN)
[74]	Fuel Consumption Estimation	Automation and Control Systems, Maintenance and Condition Monitoring, Navigational and Climate Features	Fuel Consumption [t/day]	RF, SVR	MAE: 1.2 (RF) – 2.5 (SVR) MAPE: 8% (RF) – 15% (SVR) RMSE: 1.8 (RF) – 3.1 (SVR) Fuel savings: 2%±7%
[75]	Fuel Consumption Estimation	Automation and Control Systems, Vessel – Environment Interaction, Navigational Data, Climate Features	Fuel Consumption [t/day]	ANN, SVR, PR	R <sup>2</sup> : 0.98 – MSE: 0.19 (ANN) R <sup>2</sup> : 0.49 – MSE: 1.85 (SVR) R <sup>2</sup> : 0.36 – MSE: 0.65 (PR)
[76]	Fuel Consumption and CO <sub>2</sub> Emissions Reduction	Sailing/Engine speed, Position and Fuel Consumption (Endogenous) – Wind speed/direction and Water depth/speed (Exogenous)	Fuel Consumption [Kg] and CO <sub>2</sub> Emissions [kg]	K-means Clustering PSO	Percentage Decrease: 3%
[77]	Energy Efficiency	Ship Characteristics and Voyage Speed (Endogenous) – Ice Concentration and other Environmental factors (Exogenous)	Energy Efficiency Operational Indicator (EEIO) [t/t-nm]	ANN	Accuracy error less than 5%
[78]	Fuel Consumption Estimation	Ship states (Endogenous) and Weather/Environment conditions (Exogenous)	Fuel Consumption [mt/day]	ANN, SVR, GP Lasso Regressor	RMSD: 19.5 MAE: 15.0 (ANN) RMSD: 18.7 – MAE: 13.5 (SVR) RMSD: 27.5 – MAE: 23.4 (GP) RMSD: 7.4 – MAE: 4.9 (Lasso)
[79]	Fuel Consumption Estimation	Ship's Speed (Endogenous) – Water Depth, Wind Speed, Wave, Swell and Current (Exogenous)	Fuel Consumption [mt]	ANN-MR offline ANN-MR JIT	Fuel saving: –0.43% (offline) Fuel saving: 21.24% (JIT)
[80]	Ship Speed and Engine Power Estimations	Weather Data (Exogenous) – Ship Motion and Engine States (Endogenous)	Ship Speed (SS) [kn] Engine Power (EP) [kW]	GP regressor + Ship propulsion domain knowledge	RMSE: 0.33 (SS) – 386 (EP) NRMSE: 2.23% (SS) – 2.26% (EP)

regression model which considered the current state of the vessel and the forecasted exogenous conditions. The combination of semi-empirical methods and dynamic optimization was effective in increasing ship efficiency by about 2% and reducing CO<sub>2</sub> emissions by up to 28% in the documented case studies. However, while the paper proposes route optimization with a constraint for total sailing time there is further work to be completed in terms of balancing the EEOI when the vessel is constrained by port, transport, or fleet management demands.

Wang et al. [68] employed a real-time ship control system to optimize vessel speed to the reduction of the energy demands. The authors proposed the use of Wavelet ANN to forecast the operating condition of a vessel using real-time sensor data. This model was used to obtain predictions for the water depth (Vessel – Environment interaction) and wind speed (Exogenous data) over a short sample time with six environmental parameters as input. The ship resistance and energy efficiency were then inferred using semi-empirical models common in literature. The authors applied this framework to a case study where the optimization framework used the forecasted resistance model to investigate a range of engine speeds (from 175 to 630 r/min) and selected the speed corresponding to the lowest fuel consumption. By this method, the authors proposed the fuel consumption per unit distance can be reduced by up to 19.04% in the best case. The framework presents a viable method for energy saving. However, the authors' development of the wavelet neural network was not extended to forecasts in advance of one-time step in the future. Additionally, the authors do not quantify the accuracy of their savings by additional validation methods.

Le et al. [69] apply a multi-layer perceptron framework for the prediction of vessel fuel consumption. An energy-saving framework is proposed through the prediction of fuel consumption from limited Vessel – Environment interaction data parameters including the average speed, sailing time, maximum capacity, and cargo weight. The authors suggest that this model can drive a slow steaming controller following a robust features engineering phase using the dropout procedure. The methodology is partially consistent with the literature as the authors do employ ten-fold cross-validation for the grid search model selection phase of their experiment. However, they only present the results for one split of their data with no confidence interval. The authors compare the results of the framework when applied to a containership with the fuel consumption prediction for varying semi-empirical and data-driven models with both limited and extensive features. As expected, by using an ANN and tuning the hyperparameters with a state-of-the-art approach, the accuracy of the data-driven approach supersedes the performance of the semi-empirical models by ~10% for the container vessels. The authors suggest in their conclusion that the data-driven approach is limited since it does not incorporate a priori information from the designers. Instead, they advocate that GP may be favored in the future.

Soner et al. [70] applied a tree-based method to predict the operational performance of a vessel using a dataset describing the operating conditions of a vessel. The

data included Vessel – Environment Interaction features and features for the Exogenous conditions, the authors used these parameters to predict the fuel consumption. The authors first applied a feature ranking method to determine the variables with the highest potential for information extraction. The authors compared the results from Bagging (Bag), RF, and Boostap (Boo) and were relatively comparable for the fuel consumption prediction. In particular, a Root Mean Square Error (RMSE) equal to 45.2 [L/h] for the Bag, 43.5 [L/h] for the RF, and 41.3 [L/h] for the Boo were reported. However, only one splitting of the data is presented in the paper. The authors proposed tree-based methods to develop an energy-saving framework over the use of alternative ML approaches such as the popular ANN. However, alternative algorithms in the literature performed close to those presented in this paper, and without the interval of confidence, it is difficult to make a valid comparison. Additionally, no features engineering has been performed in this paper, as is often present with a tree-based approach to improve the accuracy of regressors.

Pagoropoulos et al. [71] highlighted the importance and benefits of predictive analytics in driving energy efficiency on board, by identifying the presence of energy efficient operations. The authors focused on the electricity production from Diesel-generator sets on a group of tanker vessels by analyzing a set of endogenous and exogenous data, originating from equipment specifications about the main consumers on-board, and noon reports. The latter include a limited subset of measurements from automation and control systems, and metocean features, for over a period of two months. Based on this information, they utilized Penalized Linear Discriminant Analysis (LDA) and a SVM to solve a multi-class classification problem. Specifically, the aim was to separate between cases in which the vessel was sailing according to proper procedures, and cases in which equipment was operated inefficiently. The feasibility of the approach was demonstrated during port and sailing conditions, with an average accuracy of 98% and 94%, respectively. Although the authors did not perform systematic comparisons and a proper error estimation and model selection procedure, they underlined the potential benefits of integrating more measurements, and expanding to data streams from auto logging systems would have a very positive effect on the obtained results.

Parkes et al. [72] utilized shallow and deep feed forward ANN to estimate propulsive power demand at different environmental conditions, as a basis for more sophisticated prescriptive analytics solutions. For instance, improving weather routing and establishing power margin for new shipbuilding projects. In particular, 27 months of performance monitoring data from 3 vessels sampled every 5 min were utilized. Few endogenous and exogenous variables were available for the analysis, and these included GPS speed and speed through water, wave height, wind speed and direction, draft, trim, and heading. The authors considered multiple ANN architectures, focusing primarily on the number of layers and neurons in each layer, for a total of 16 architectures. Through the proper error estimation and model selection procedures, the actual propulsive power could be

predicted with a Mean Absolute Percentage Error (MAPE) of 10% on average for all architectures, with the best ANN architecture resulting in a MAPE of 8%. The authors concluded that the quality of the predictions is satisfactory, underlying that the extent to which the ANN captures the underlying physics of ship performance is still an open issue.

Xu et al. [73] demonstrated the application of predictive analytics for wear fault diagnosis of marine Diesel engines. The authors developed a multi-model fusion system, considering the predictions of an Evidential Reasoning (ER) model, a Belief Rule-Based (BRB) inference model, and an ANN, as pieces of evidence to be fused in a decision level, for the classification of 8 wear-and-tear faults. The diagnostic framework also included a novel methodology to determine the reliability of the predictions of each model, and the use of a genetic algorithm to assign relative importance weights to the predictions of each model. Utilizing a total of 150 samples acquired from the control and condition monitoring systems of marine diesel engines in operation, and through proper error estimation procedures, the authors showed that the accuracy of their diagnostic framework could correctly classify faults with an accuracy ranging between 93–100%, an increase of approximately 8%, compared to the smallest accuracy given by the single models.

Yan et al. [74] demonstrated the advantages of adopting prescriptive analytics in the shipping industry. They proposed a two-stage sailing speed optimization framework for a dry bulk vessel, with the aim of minimizing the vessel's fuel consumption during a voyage subject to the required arrival time at the port of destination. Decision Tree (DT), ANN, RF, SVM, and Lasso Regressors were developed and validated to estimate the vessel's fuel consumption, considering sailing speed, total cargo weight, and climate features, utilizing information available in the noon reports for a time span of 2 years. Although the authors have not reported detailed results regarding model selection and error estimation, they obtained a Mean Absolute Percentage Error (MAPE) of 7.91% in fuel consumption estimation. Subsequently, they formulated a mixed-integer linear programming optimization problem that uses the predictions of the RF, to derive the optimal vessel speed throughout the voyage. They compared the obtained solutions with real historical data on two 8-day voyages and concluded that fuel savings of 2%–7% can be achieved. Furthermore, it was argued that developing more sophisticated optimization models would help to yield more practical management strategies. The authors further underlined the importance of accurately estimating the environmental conditions as they heavily influence the quality of the obtained solutions, and highlighted the importance of obtaining additional endogenous and exogenous features to develop and benchmark their RF regression algorithm.

Jeon et al. [75] developed an ANN for the prediction of the fuel consumption of a marine Diesel engine, which they also compared with other ML models, including SVR and Polynomial Regression (PR). The authors utilized a limited set of both endogenous



and exogenous data for the fuel consumption estimation, including information from the navigational and cargo monitoring systems, climate features, and engine speed and power output. However, they do not provide all details regarding the data collection process. A systematic variation on the number of layers and the number of neurons in each layer was performed, ranging between 0–2 and 1–7 respectively, with three activation functions. The performance of all architectures was benchmarked on a single subset of data, without detailing the error estimation and model selection procedures. It was concluded that through proper calibration of the parameters of an ANN, highly accurate results can be obtained, that surpass the performance of the SVR or PR, although the variance of the results was high. More specifically, the best performing ANN architecture is characterized by an MSE of 0.19 and an  $R^2$  of 0.98, whereas for SVR and PR the results were not satisfactory, and possibly a more thorough model selection could have helped achieve better results.

Yan et al. [76] apply data analytics techniques to reduce vessel transportation's CO<sub>2</sub> emissions that correspond to approximately 2.7% of the global releases. Inland ships are included in this survey and the authors focus on proposing an optimization workflow, based on big data analysis, to minimize this issue. They use a K-means clustering algorithm to classify each segment of the ship's route based on exogenous factors. Then a ship energy efficiency model considering both exogenous and endogenous factors is set up. Finally, a swarm intelligence algorithm based on iterative processes, the Particle Swarm Optimization algorithm (PSO), is adopted to solve the non-linear optimization problem concerning the optimal engine speeds under different environmental conditions. The authors test their approach on the Yangtze river case study and they achieve a reduction in fuel consumption and CO<sub>2</sub> emissions over an entire trip equal to 3% and 2.38% during the dry and rain seasons, respectively. The presented optimization method can be extended to include others or even more influencing factors, such as ship parameters, route characteristics, port operation, and transport demand.

Zhang et al. [77] focused on transportation in Arctic waters which is becoming very attractive nowadays thanks to shorter and faster sea routes connecting main continents. Nonetheless, the fragility and sensitivity of the Arctic environment require a feasibility analysis and ship energy efficiency improvements. For this reason, the authors propose a methodology based on three steps. Firstly, they perform a Pearson correlation analysis to find the energy efficiency's most influencing factors. Then, they feed a prediction model based on an ANN with previously detected variables to find the best ship speed optimization strategy. The prediction model achieves a 0.98 Fitting Rate on the test set that corresponds to an accuracy error less than 5%. Finally, the improved Ant Colony Algorithm (ACA) is adopted to solve the optimum energy efficiency route planning problem. The authors test this approach on Yong Sheng ship's collected data, and the results point out the importance of considering also the energy efficiency during ship's route planning instead of the distance uniquely. Further improvements can be achieved

enriching network's inputs with a greater number of parameters. Moreover, this analysis would benefit from a validation on other case study and, finally, since the concrete danger of navigating in Arctic waters, a risk analysis needs to be included in this optimization methodology.

Wang et al. [78] focused on the maritime fuel consumption and optimization. Their aim was to develop a novel predictive model able to estimate the fuel consumption of a specific ship as a function of the ship's state and surrounding environments. They clearly pointed out the several numbers of feature variables affecting the fuel consumption characterized by colinearities which make a traditional multiple linear regression method unable to correctly predict the fuel consumption. To deal with this problem, the authors employ the Lasso Regression algorithm to implement the variable selection of feature variables and to build an accurate ship fuel consumption prediction model. The authors show that the latter methodology can successfully select 20 features out of more than 30. Moreover, to demonstrate the performance of the proposed predictive model, they compare it with typical regression methods, such as ANN, Support Vector Regression (SVR), and GP. They decide to measure the performances of the latter models in terms of MAE and Root Mean Square Deviation (RMSD). The proposed Lasso-based model outperforms the other three models on the same test set with a RMSD of 7.4 [mt/day] and a MAE of 4.9 [mt/day]. Thus, it is able to fit accurately the real values most of situations. To further improve the model's performance a greater data collection of different voyages or ships is needed.

Farag and Ölçer [79] developed a combined ANN and Multi-Regression (MR) model able to predict ship fuel consumption under different sea environment conditions in a real-time varying scenario. The authors select the Brake Power feature variable as the target of their ANN model, after performing a correlation analysis on the available dataset. They demonstrate how the proposed combined prediction model outperforms a single ANN-based model. Moreover, they employ the validated performance model to assess the potential savings of real-time, or rather Just-In-Time (JIT), measures and they successfully achieve a 24.24% of heavy fuel oil and 328.5 ton of CO<sub>2</sub> emissions savings. Further improvements of the actual model can be accomplished expanding the available dataset and considering other modeling approaches.

Yoo and Kim [80] analyze the powering performance of full-scale ships under varying environmental and operating conditions to face the increased concerns about environmental pollution and global warming. Since classical ML algorithms are often susceptible to statistical overfitting in dealing with such data sources with many influencing factors, they propose to incorporate domain knowledge of ship propulsion into the design of two regression models able to predict the optimal Ship Speed (SS) and Engine Power (EP). They choose to employ GP regression model thank to its effectiveness in describing complicated nonlinear models. After identifying domain influencing factors through graphical models' exploitation, the authors employ the regularization

scheme to integrate such domain knowledge. For the training dataset, the GP model without regularization outperforms the GP with regularization in terms of RMSE and Normalized Root-Mean-Square Error (NRMSE). For the test dataset, on the contrary, the GP model with regularization achieves the best results. This confirms that the regularization technique prevents overfitting and leads to better regression performance. The authors estimate that the travel time and total fuel consumption can be predicted over any path segment and that the proposed method can be directly applied to optimal weather routing.

## 7.4. Case studies

In the following, the authors provide an in-depth analysis of three relevant case studies:

1. Data-driven condition monitoring of a marine dual fuel engine;
2. Data-driven digital twin to estimate the marine fouling status;
3. Trim optimization to reduce fuel consumption.

### 7.4.1 Data-driven condition monitoring of a marine dual fuel engine

Maritime transportation accounts for around 80% of the world freight movements, remarkably contributing to the global environmental footprint. Dual fuel engines, running on both gaseous and liquid fuels, represent a viable way toward the reduction of emissions at the cost of additional complexity in monitoring activities. While various traditional approaches to monitoring exist, data-driven methods represent the frontier in research and in maritime industrial applications. Data-driven monitoring methods usually require a large amount of labeled data, i.e., sensor measurements plus the associated engine status usually annotated by human operators, which are costly and seldomly available in the wild. Unlabeled samples, instead, are commonly, cheaply, and readily available. The enabling technology for data-driven methods is the availability of a network of sensors and an automation system able to capture and store the associated stream of data. In [81], authors design and propose multiple alternatives toward the weakly supervised marine dual fuel engines data-driven monitoring. To this aim, the authors developed a Digital Twin of the dual fuel engine and novelty detection algorithms. Results on data generated from a real-data validated simulator of a marine dual fuel engine demonstrated that the proposed weakly supervised monitoring approaches lead to a negligible loss in accuracy compared to costly and often unfeasible fully supervised ones supporting the validity of the proposal for its application in the wild. The main outcome of the paper was to guide researches and practitioners for the selection of the best data-driven dual fuel engine monitoring method according to the available data about the vessel.

### 7.4.1.1 The approach

The approach employed in [81] includes the following three steps:

1. Fully Supervised Performance estimation;
2. Fully Supervised Health Status Estimation;
3. Weakly Supervised Health Status Estimation.

The Fully Supervised Performance Estimation step includes the design of a Digital Twin, exploiting state-of-the-art supervised data-driven methods for enabling the prediction of the engine performance and emissions parameters based on the control variables (e.g. engine load and engine speed), in healthy engine conditions. This step actually does not employ labeled data; instead it employs the acquired data from engine operation under healthy conditions. The Fully Supervised Health Status Estimation step focuses on developing models capable of classifying the status of the engines as healthy or faulty, and it is accomplished by employing two approaches. The first one employs the Digital Twin developed in the first step to estimate the deviation (drift) of the parameters of the actual engine operation (based on the acquired data) from the respective Digital Twin predicted parameters. The second one exploits state-of-the-art supervised data-driven methods to classify the status of the investigated engine based on the control and performance parameters. This step requires labeled data of the engine under healthy and faulty conditions. The Weakly Supervised Health Status Estimation step focuses on reducing the amount of labeled data required to build the models developed in the second step by employing two approaches. The first one focuses on the estimation of the engine performance parameters variation from the respective parameters calculated by employing the Digital Twin by employing a limited amount of labeled data for tuning the drift detection model. The second one, instead, exploits state-of-the-art unsupervised data-driven methods to detect abnormal conditions (anomalies) of the investigated engine by employing as input the considered control and performance parameters. The weakly supervised health status estimation step employs the models trained just with data acquired under the engine healthy conditions from the engine monitoring system. These models are subsequently fine tuned with a very small amount of labeled data. Fig. 7.4 depicts the authors proposal with a simple graphical representation.

### 7.4.1.2 Data description

Since datasets corresponding to the investigated marine DF engine under faulty conditions were not available to the authors (for the reason described in the original manuscript) this study employed a high fidelity physical model developed and validated in previous authors' studies [82] to generate the data. For this purpose, multiple simulation runs, corresponding to different scenarios, were performed collecting engine control and performance parameters to generate this dataset.

Two of these datasets correspond to the engine operation at healthy conditions in both the diesel mode ( $DB_{Healthy}^{DM}$ ) and the gas mode ( $DB_{Healthy}^{GM}$ ). These datasets were

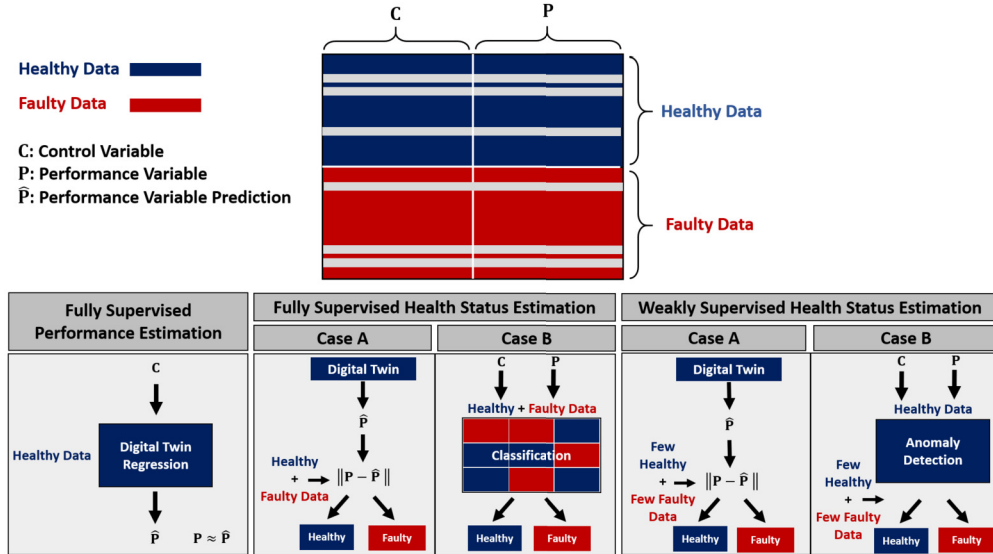


Figure 7.4 Case Study 1 – The approach.

referred as the “healthy status” datasets. Moreover, a third dataset ( $DB_{Faulty}^{DM}$ ) is created that contains the control and the performance variables corresponding to the degraded conditions (described in the original manuscript) and the engine diesel mode. This dataset was referred as the “faulty status” dataset.

The control (input –  $C$ ) and performance (output –  $P$ ) variables collected from the physical model and included in the datasets are listed in Table 7.4.

### 7.4.1.3 Results

In this section, we summarize the results of [81] in applying the approach summarized in Section 7.4.1.1 over the data summarized in Section 7.4.1.2.

For what concerns the Fully Supervised Performance estimation scenario, Table 7.5 reports the best error (measured with the MAPE) over three different models (RF, KM, and ANN) for both diesel and gas modes. As reported in Table 7.5, the MAPE of the best performing model (KM) is always less than 4% (and in most case less than 2%) for all the performance variables, in the considered modalities.

For what concerns the Fully Supervised Health Status Estimation scenario, three different models (RF, KM and ANN) have been applied again to distinguish between two possible engine’s statuses: +1 for faulty condition and -1 for healthy condition. Table 7.6a reports the misclassification errors of the best of the tree model for the two different approaches under analysis (*Direct* and *Digital Twin*). As one can observe from Table 7.6, the *Direct* approach outperforms the *Digital Twin* one in terms of misclassi-

**Table 7.4** Case Study 1 – The data.

Type	Name	Unit
$C_1$	Ambient Temperature	[K]
$C_2$	Air Cooler Temperature	[K]
$C_3$	Gas Valve Unit Gas Pressure	[bar]
$C_4$	Heating Value of Diesel	[MJ/kg]
$C_5$	Heating Value of Gas	[MJ/kg]
$C_6$	Engine Load	[kW]
$C_7$	Engine Speed	[rpm]
$P_1$	Fuel Rack Position (Diesel Mode)	
$P_2$	Main Gas Pressure	[bar]
$P_3$	Max Cylinder Pressure	[bar]
$P_4$	Charged Air Pressure (Inlet section)	[bar]
$P_5$	Exhaust Gas Temperature Turbocharger – inlet	[K]
$P_6$	Turbocharger Speed	[rpm]
$P_7$	Waste Gate Opening	[mm <sup>2</sup> ]
$P_8$	Nitrogen Oxide	[g/kWh]
$P_9$	Carbon dioxide	[g/kWh]
$P_{10}$	BSFC	[g/kWh]

**Table 7.5** Case Study 1 – Fully Supervised Performance estimation – Error, measured with the MAPE, of the best model (among RF, KM, and ANN) in predicting the performance variables for both diesel and gas mode.

Engine $P$	Diesel mode	Gas mode
$P_1$	0.31±0.02	0.00±0.00
$P_2$	0.91±0.05	2.67±0.15
$P_3$	0.22±0.01	3.56±0.21
$P_4$	0.49±0.03	2.97±0.17
$P_5$	1.37±0.07	4.20±0.30
$P_6$	1.17±0.06	2.38±0.13
$P_7$	–	3.06±0.18
$P_8$	1.77±0.12	1.46±0.08
$P_9$	0.55±0.03	1.13±0.07
$P_{10}$	1.71±0.10	1.39±0.08

fication errors for all the exploited models. In order to better represent the quality of the developed model, Tables 7.6b and 7.6c report the confusion matrices for the best models and for both the adopted approaches. By observing these confusion matrices it is possible to note that the misclassification errors are well distributed and models do not tend to predict more false positive than false negative. Also in this case the quality of

**Table 7.6** Case Study 1 – Fully Supervised Health Status Estimation - Error (percentage of misclassification error) and confusion matrices (positive are faulty and negative are healthy) of the best model (among RF, KM, and ANN) with a Direct Approach or with a Digital Twin in predicting the health status conditions of the engine in diesel mode.

(a) Misclassification error %

Engine	Diesel mode	
Approach	Direct	Digital twin
	1.7±0.1	4.3±0.2

(b) Confusion matrix with KM for diesel mode using a Direct Approach.

		Actual	
Prediction	TP	98.2±0.1	FP 1.8±0.1
	FN	1.6±0.1	TN 98.4±0.1

(c) Confusion matrix with KM for diesel mode using a Digital Twin.

		Actual	
Prediction	TP	96.2±0.2	FP 3.8±0.2
	FN	4.8±0.3	TN 95.2±0.3

the developed models is surely up to a level which is acceptable for their use in the wild with misclassification below the 5% (and in most case less than 3%). It is worth noting that switching from the Fully Supervised Performance estimation scenario to the Fully Supervised Health Status Estimation scenario does not compromise the ability to make accurate predictions.

For what concerns the Weakly Supervised Health Status Estimation scenario Table 7.7a reports the misclassification error percentages and the related confusion matrices of the best model. Diesel engine mode data and the Direct and Digital Twin approaches are exploited in this context, in accordance with previous scenario's experiments. In this case, the experimental setup is slightly different, by drastically reducing the need for labeled data. For the Direct approach, One Class Support Vector Machines (OCSVM) and the Global KNN (GKNN) models are taken into consideration and compared against the ones based on the Digital Twin (built with RF, KM, and ANN). In this setting, the validation set size (which is the only actual labeled one) is a parameter to be considered as another degree of freedom (which must be kept as small as possible to be able to apply this methodology in the wild). Observing Table 7.7a, the Direct approach outperforms again the Digital Twin-based one. For the Direct approach, OCSVM model outperforms the GKNN one. Meanwhile, for the Digital Twin-based

**Table 7.7** Case Study 1 – Weakly Supervised Health Status Estimation - Error (percentage of misclassification error) and confusion matrices (positive are faulty and negative are healthy) of the best model with a Direct Approach (among OCSVM and GKNN) or with a Digital Twin (among RF, KM, and ANN) in predicting the health status conditions of the engine in diesel mode.

(a) Misclassification error %

Engine	DM	
	Direct	Digital twin
$n_v = 10$	$2.3 \pm 0.1$	$4.9 \pm 0.3$
$n_v = 20$	$2.2 \pm 0.1$	$4.7 \pm 0.2$
$n_v = 40$	$2.1 \pm 0.1$	$4.4 \pm 0.2$

(b) Confusion matrix with OCSVM for diesel mode using a Direct Approach ( $n_v = 10$ ).

		Actual		
Prediction	TP	$97.6 \pm 0.1$	FP	$2.4 \pm 0.1$
	FN	$2.2 \pm 0.1$	TN	$97.8 \pm 0.1$

(c) Confusion matrix with KM for diesel mode using a Digital Twin ( $n_v = 10$ ).

		Actual		
Prediction	TP	$95.5 \pm 0.2$	FP	$4.5 \pm 0.2$
	FN	$5.3 \pm 0.3$	TN	$94.7 \pm 0.3$

approach, KM model confirms to outperform both RF and NN models. Increasing the validation set size (the amount of labeled data) increases also the performance. For the Direct approach the improvement is not so relevant, while for Digital Twin approach is more relevant. This confirms the higher ability of the Direct approach to deliver high performance with limited number of labeled samples. Tables 7.7b and 7.7c report the confusion matrices of the best models under analysis for a validation set size equal to 10. Also in this case, observing these confusion matrices it is possible to note that the misclassification errors are well distributed and models do not tend to predict more false positive than false negative. Also in this case the quality of the developed models is surely up to a level which is acceptable for their use in the wild with misclassification below the 5% (and in most case less than 3%). Note then that, switching to this last scenario does not compromise the ability to make accurate predictions (the decrease in performance is less than 1%).



#### **7.4.1.4 Conclusions**

In this paper, authors of [81] focus on data-driven monitoring models to be employed in the wild. Unfortunately data-driven methods often require a large amount of labeled samples which are rarely available. For this reason, authors design and propose multiple alternatives toward the weakly supervised marine dual fuel engines data-driven monitoring. Results on data generated from a real-data validated simulator of a marine dual fuel engine demonstrate that the proposed weakly supervised monitoring approaches lead to a negligible loss in accuracy compared with costly and often unfeasible fully supervised ones supporting the validity of the proposal for its application in the wild. In particular, in the Fully Supervised Performance estimation scenario, the error of the data-driven model is always less than 4% (and in most cases less than 2%) for all the performance variables in the considered modalities. This result is surely up to a level which is acceptable for the utilization of data-driven models for dual fuel engine performance estimation. Considering the Fully Supervised Health Status scenario, the error of the fault detection models is always below the 5%, and in most cases less than 2% which is again suitable for real operational environment, but unfortunately it requires a number of labeled samples which is not realistic to obtain in the wild. Finally, in the Weakly Supervised Health Status Estimation scenario, we fill this gap by remarkably decreasing the amount of labeled samples necessary to train the model whilst obtaining an error below the 5% (and in most cases less than 3%) and not compromising the ability to make accurate predictions (the decrease in performance is less than 1%) for the use of this model in real operational conditions.

#### **7.4.2 Data-driven digital twin to estimate the marine fouling status**

Shipping is responsible for approximately the 90% of world trade leading to significant impacts on the environment. As a consequence, a crucial issue for the maritime industry is to develop technologies able to increase the ship efficiency, by reducing fuel consumption and unnecessary maintenance operations. For example, the marine fouling phenomenon has a deep impact, since to prevent or reduce its growth which affects the ship consumption, costly drydockings for cleaning the hull and the propeller are needed and must be scheduled based on a speed loss estimation. In [83] a data driven Digital Twin of the ship is built, leveraging on the large amount of information collected from the on-board sensors, and is used for estimating the speed loss due to marine fouling. A thorough comparison between the proposed method and ISO 19030, which is the de-facto standard for dealing with this task, is carried out on real-world data coming from two Handymax chemical/product tankers. Results clearly show the effectiveness of the proposal and its better speedloss prediction accuracy with respect to the ISO 19030, thus allowing reducing the fuel consumption due to fouling.

**Table 7.8** Case Study 2 – Main features of V1 and V2 case studies.

Ship feature	V1		V2	
	Value	Unit	Value	Unit
Deadweight	46764	[t]	46067	[t]
Design speed	15	[knots]	15.5	[knots]
Draft (summer SW)	12.18	[m]	12.2	[m]
Length between perpendicular	176.75	[m]	176.83	[m]
Breadth moulded	32.18	[m]	32.20	[m]
Main engines installed power	3840×2	[kW]	8200	[kW]
Auxiliary engines installed power	682×2	[kW]	1176×3	[kW]
Shaft generator power	3200	[kW]		
Exhaust boilers steam generator	750×2	[kg/h]	1130	[kg/h]
Auxiliary boilers steam generator	14000×2	[kg/h]	14000×2	[kg/h]
Fuel consumption	34.7	[mt/day]	31.8	[mt/day]

#### 7.4.2.1 The approach

Inspired by the ISO 19030 [84] and supported by the evidence that data-driven models can be much more accurate and effective than the physical ones, in this work the authors proposed a data-driven model for predicting the vessel's speed, able to act as a “Digital Twin” [85] of the ship herself. The Digital Twin can be used to compute the deviation between the predicted performance and the actual one, namely the speed loss [86]. The authors showed that the average drift in time of the speed loss can be exploited to accurately and effectively estimate the effects of the marine fouling on the ship performance, and thus program a more efficient hull and propeller cleaning scheduling. To this aim, authors of [83] propose a two-phase approach:

1. a Digital Twin based on a data-driven model, is built using a deep neural networks and the data described in Section 7.4.2.2. The model exploits data collected during a suitable period of time when the marine fouling is not present and for a period long enough to observe the ship in different operational and environmental conditions (e.g., one can start the data collection just after the launch of the ship or its hull and propeller cleaning and stop after one or two months of operations);
2. the data-driven model is applied on a second set of data and the speed loss is computed. Subsequently, the drift in the average behavior of the speed loss between two maintenance operations is studied, together with changes in its distribution using robust regression and statistical nonparametric test.

#### 7.4.2.2 Data description

This section presents the two Handymax chemical/product tankers exploited (see Table 7.8 for the main features of the vessels) and the available data.

**Table 7.9** Case Study 2 – Data collected from logging system of the two vessels.

Variable name	Unit	Variable name	Unit
Timestamp	[t]	Sea depth	[m]
Latitude	[°]	Seawater temperature	[°C]
Longitude	[°]	CPP set point	[°]
Main engines fuel consumption	[kg/h]	CPP feedback	[°]
Auxiliary engines power output	[kg/h]	Fuel density	[kg/m <sup>3</sup> ]
Shaft generator power	[kg/h]	Fuel temperature	[°C]
Propeller shaft power	[kW]	Ambient pressure	[bar]
Propeller speed	[rpm]	Humidity	[%]
Ship draft (fore)	[m]	Dew point temperature	[°C]
Ship draft (aft)	[m]	Shaft torque	[kN m]
Draft port	[m]	Rudder angle	[°]
Draft starboard	[m]	Acceleration x direction	[m/s <sup>2</sup> ]
Relative wind speed	[m/s]	Acceleration y direction	[m/s <sup>2</sup> ]
Relative wind direction	[°]	Acceleration z direction	[m/s <sup>2</sup> ]
GPS heading	[°]	Roll	[°]
Speed over ground	[knots]	Pitch	[°]
Speed through water	[knots]	Yaw	[°]

The two vessels are equipped with the same data logging system which is used by the company for both on board monitoring and land-based performance control. Table 7.9 summarizes the available measurements from the continuous monitoring system. The original frequency of data acquisition by the monitoring system is equal to 1 point every 15 seconds.

The available data of the two vessels have been collected in the time slots for V1, between the 21/03/2012 17:45:00 and the 03/10/2014 14:15:00, and for V2 between 01/05/2014 00:15:00 and 26/08/2016 14:15:00.

At last, Table 7.10 reports the recorded relevant maintenance events of the two vessels.

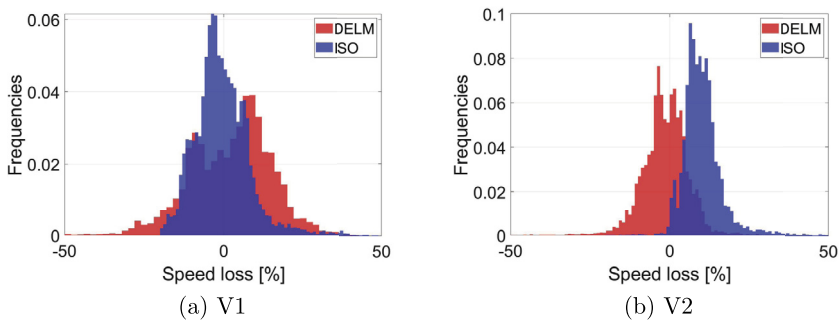
### 7.4.2.3 Results

As results will show the data-driven model proposed by the authors of [83] allows the identification of clear drift in the performance of the vessel compared to the ISO 19030 procedure.

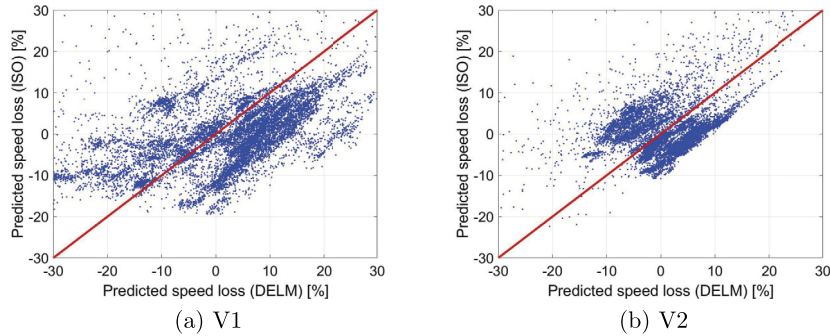
In Fig. 7.5 the histograms of the percentage speed losses computed with ISO 19030 of the percentage speed losses for V1 and V2 are reported. The variance of the distribution of the percentage speed losses is larger for the data driven model with respect to the one of the ISO 19030. This is caused by the fact that the ISO 19030 filters out a large amount of data points, only keeping those for which the application of the method

**Table 7.10** Case Study 2 – Maintenance events for V1 and V2.

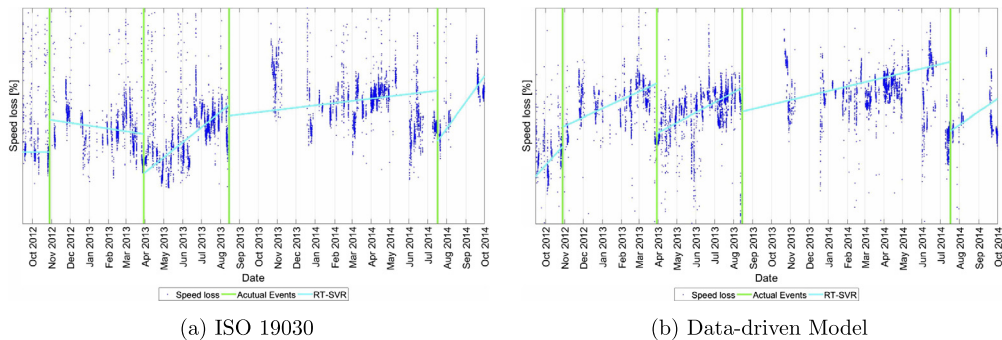
Date	Event
<b>V1</b>	
21/03/2012	Vessel delivery
29/10/2012	Propeller cleaning
30/03/2013	Hull cleaning
01/08/2013	Loss of the LOG speed measurement
17/07/2014	Change from fixed-speed to variable-speed operations
<b>V2</b>	
19/04/2014	Propeller polishing
20/12/2014	Hull cleaning
28/08/2015	Hull cleaning and Propeller polishing
28/11/2015	Dry-docking

**Figure 7.5** Case Study 2 – Histograms of data-driven model and ISO 19030 Estimated Percentage Speed Loss.

is more reliable. On the other hand, the data-driven model exploits all the available data points corresponding to a larger variety of operational conditions. Moreover, the average of the distribution of the speed loss is not always centered on a positive value, due to the fact that the data used for training the data-driven model and the parameters used for the ISO 19030 do not correspond to a perfect clean state, as it would be required for creating a perfect digital twin (as shown later, this problem does not affect the quality of the final results). Finally, the results obtained by the two models are, at least qualitatively, in an overall good agreement. In order to better quantify the agreement between the data-driven model and the ISO 19030 models, Fig. 7.6 reports the scatter-plot of the data-driven model and the ISO 19030 estimated percentage speed loss for V1 and V2. From Fig. 7.6 it is possible to observe that the speed loss predicted by the data-driven model and the ISO 19030 methods are positively correlated (particularly for V2), thus demonstrating that the prediction achievable by the proposed data-driven model is consistent with the state-of-the-art.

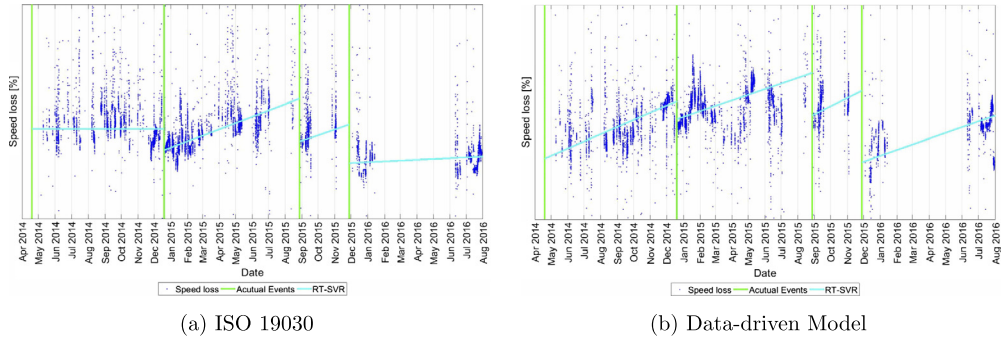


**Figure 7.6** Case Study 2 – Scatterplot of the data-driven model and the ISO 19030 Estimated Speed Loss Percentages.

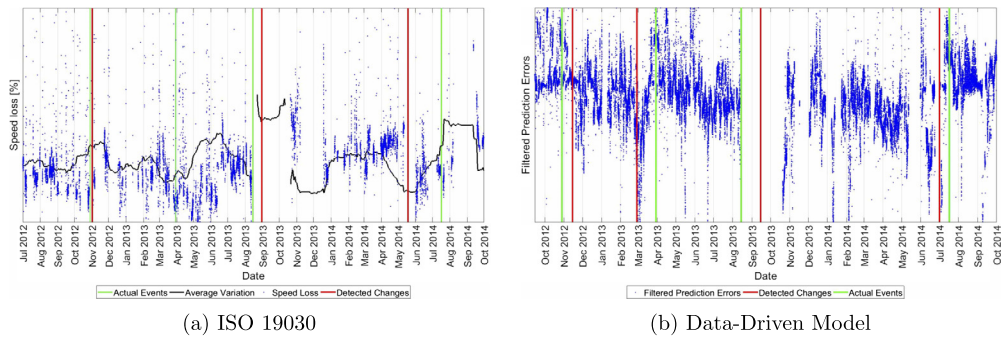


**Figure 7.7** Case Study 2 – Robust Regression Analysis on the Speed Loss Percentages between two consecutive Hull and Propeller Cleaning Events for V1.

Then authors of [83] report the analysis of the drift in data-driven model and ISO 19030 estimated percentage speed loss between two consecutive hull and propeller cleaning events, carried out with the robust regression analysis. Figs. 7.7 and 7.8 report the results for V1 and V2 respectively. Those results clearly show the higher level of reliability of the prediction achieved by the data-driven model method against the ISO 19030 one. In both vessels, the linear trend for the speed loss calculated by the ISO 19030 method shows large variations between different maintenance intervals. In addition, in some intervals between two consecutive hull and propeller cleaning operations, the trend in the estimated percentage speed loss using the ISO 19030 method is negative. These results do not agree with the physical basis of the fouling phenomenon, and suggest that, in the case presented in this paper, the application of the ISO can lead to inaccurate results. On the contrary, as far as data-driven model is concerned, Figs. 7.7 and 7.8 clearly show trends that are always physically plausible. Model drift behavior be-



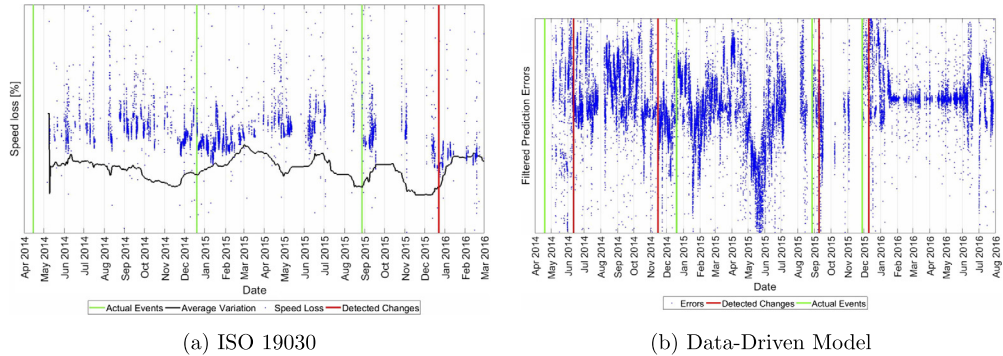
**Figure 7.8** Case Study 2 – Robust Regression Analysis on the Speed Loss Percentages between two consecutive Hull and Propeller Cleaning Events for V2.



**Figure 7.9** Case Study 2 – Changes in time of the distribution of the percentage of speed loss estimated with the ISO 19030 and the Data-Driven Model for V1.

tween different cleaning intervals is now consistent with the one characterizing a ship that operates in conditions, on average, similar over time.

Finally authors of [83] report the analysis of the changes in time of data-driven model and ISO 19030 estimated percentage speed loss distributions, carried out with the nonparametric statistical test of Kolmogorov–Smirnov. Figs. 7.9 and 7.10 report the results for V1 and V2 respectively. Those figures testify the higher level of reliability of the data-driven model method against the ISO 19030 one. In both vessels there are no statistically meaningful changes in the distribution of the speed losses estimated with the ISO 19030, and just in a few cases the Kolmogorov–Smirnov test detects a change in correspondence to an actual hull and propeller cleaning event (see Table 7.10). On the contrary, when the same method is applied to the speed losses estimated with the data-driven model, the Kolmogorov–Smirnov test detects all the changes in correspondence or close to an actual hull and propeller cleaning event.



**Figure 7.10** Case Study 2 – Changes in time of the distribution of the percentage of speed loss estimated with the ISO 19030 and the Data-Driven Model for V2.

#### 7.4.2.4 Conclusions

In this work authors of [83] focus on the problem of estimating the speed loss caused by the effect of fouling on the ship hull and propeller. For this purpose, a novel data-driven method in combination with advanced statistical methods is proposed. Results on real-world data coming from two Handymax chemical/product tankers show the effectiveness of the proposal and its better prediction accuracy and reliability, with respect to the ISO 19030 de facto standard. The proposal shown to be both more accurate in predicting the loss of performance over time, between cleaning intervals, and in automatically detecting maintenance events.

### 7.4.3 Trim optimization to reduce fuel consumption

The authors of [87] investigated the problems of predicting the fuel consumption and of providing the optimal value for the trim of a vessel in real operations based on data measured by the onboard automation systems. To this aim, the authors developed and compared three different approaches for the prediction of the fuel consumption: White, Black and Gray Box Models. White Box Models (WBM) are based on the knowledge of the physical underlying processes. Black Box Models (BBMs) build upon statistical inference procedures based on the historical data collection. Finally, the authors propose two different Gray Box Model (GBM) which are able to exploit both mechanistic knowledge of the underlying physical principles and available measurements. Based on these predictive models of the fuel consumption a new strategy for the optimization of the trim of a vessel has been developed and proposed. Results on real world operational data show that the BBM is able to remarkably improve a state-of-the-art WBM, while the GBM is able to encapsulate the a-priori knowledge of the WBM into the BBM so to achieve the same performance of the latter but requiring less historical data. Further-

more, the authors proved that the GBM can be used as an effective tool for optimizing the trim of a vessel in real operational conditions.

#### **7.4.3.1 The approach**

Trim optimization has been extensively discussed in the past. It is well known, from hydrodynamics principles, that the trim of the vessel can significantly influence its fuel consumption. Previous work in scientific literature related to trim optimization has focused on three main alternative strategies: WBM and BBMs. WBMs describe the behavior of the ship resistance, propeller characteristics and engine performances based on governing physical laws and taking into account their mutual interactions. The higher the detail in the modeling of the physical equations which describe the different phenomena, the higher the expected accuracy of the results and the computational time required for the simulation. WBMs are generally rather tolerant to extrapolation and do not require extensive amount of operational measurements; on the other hand, when employing models that are computationally fast enough to be used for online optimization, the expected accuracy in the prediction of operational variables is relatively low. In addition, the construction of the model is a process that requires competence in the field, and availability of technical details which are often not easy to get access to. Differently from WBMs, BBMs make use of statistical inference procedures based on historical data collection. These methods do not require any a-priori knowledge of the physical system and allow exploiting even measurements whose role might be important for the calculation of the predicted variables but might not be captured by simple physical models. On the other hand, the model resulting from a black-box approach is not supported by any physical interpretation, and a significant amount of data (both in terms of number of different measured variables and of length of the time series) are required for building reliable models. GBMs have been proposed as a way to combine the advantage of WBMs and BBMs. According to the GBMs principles, an existing WBM is improved using data-driven techniques, either in order to calculate uncertain parameters or by adding a black-box component to the model output. GBMs allow exploiting both the mechanistic knowledge of the underlying physical principles and available measurements. The proposed models are more accurate than WBMs with similar computational time requirements, and require a smaller amount of historical data when compared to a pure BBMs. The aim of [87] was to propose the application of a GBMs to the prediction of ship fuel consumption which can be used as a tool for online trim optimization. In this framework the authors exploit data-driven methods (kernel methods and ensemble techniques) to improve an effective but simplified physical model [88] of the propulsion plant.

#### **7.4.3.2 Data description**

The data exploited in this work are the ones of Case 7.4.2 described in Section 7.4.2.2.



**Table 7.11** Case Study 3 – Indexes of performance of the WBM in predicting the Shaft Power and Fuel Consumption.

Shaft Power					
MAE [KW]	MAPE [%]	MSE [KW <sup>2</sup> ]	NMSE	REP [%]	PPMCC
7.69e+02	17.85	1.00e+06	1.13	23.59	0.65
Fuel Consumption					
MAE [ $\frac{g}{KWh}$ ]	MAPE [%]	MSE [ $\frac{g^2}{KW^2h^2}$ ]	NMSE	REP [%]	PPMCC
5.14e-02	20.95	3.94e-03	1.98	25.40	0.63

**Table 7.12** Case Study 3 – Indexes of performance of the best BBM in predicting the Shaft Power and Fuel Consumption.

Shaft Power					
MAE [KW]	MAPE [%]	MSE [KW <sup>2</sup> ]	NMSE	REP [%]	PPMCC
7.67e+01	1.90	2.47e+04	0.03	3.18	0.99
Fuel Consumption					
MAE [ $\frac{g}{KWh}$ ]	MAPE [%]	MSE [ $\frac{g^2}{KW^2h^2}$ ]	NMSE	REP [%]	PPMCC
4.62e-03	1.95	1.10e-04	0.06	3.56	0.96

### 7.4.3.3 Results

In this section we summarize the results of [87].

First we report the performance of the WBM in Table 7.11. Results show that the WBM does not show sufficient accuracy when compared with operational measurements. The inability of the model to take into account the influence of the sea state (i.e., wind and waves) on the required propulsion power is considered to be the largest source of error for this model.

Then we report the results of the best BBM (among the one tested in [87]) in Table 7.12. Note that the BBM remarkably outperform the WBM since they are able to take into account all the available information measured by the on board sensors. Among the different BBMs, the one based on ensemble methods shows the most promising results.

Then we report the results of the best GBM (among the one tested in [87]) in Table 7.13. Note that the GBMs outperform the both BBMs and GBM since they are able to take into account both the physical knowledge about the system and all the available information measured by the on board sensors. Among the different GBMs, the one based on kernel methods seems to be the best performing ones.

Observing the results on feature ranking in [87] it is also possible to note how the BBMs and GBMs actually learn meaningful information from the data which is physically plausible diminishing the doubts about the presence of spurious correlations.

**Table 7.13** Case Study 3 – Indexes of performance of the best GBM in predicting the Shaft Power and Fuel Consumption.

Shaft Power					
MAE [KW]	MAPE [%]	MSE [KW <sup>2</sup> ]	NMSE	REP [%]	PPMCC
3.18e+01	0.79	1.06e+04	0.01	1.35	0.99
Fuel Consumption					
MAE [ $\frac{g}{KWh}$ ]	MAPE [%]	MSE [ $\frac{g^2}{KW^2h^2}$ ]	NMSE	REP [%]	PPMCC
1.97e-03	0.83	4.71e-05	0.02	1.55	0.97

**Table 7.14** Case Study 3 – Fuel Consumption percentage reduction with the trim Optimization technique.

$\delta$	% reduction
0%	0.52 ± 0.12
1%	1.45 ± 0.32
2%	1.72 ± 0.51
5%	2.22 ± 0.67
10%	2.30 ± 0.64

Finally, Table 7.14 reports the Fuel Consumption percentage reduction with the trim Optimization technique proposed in [87]. Then result is reported varying the trim ( $\delta$ ) namely how much (in %) we are willing to deviate from the trim selected by the operators. As expected, the optimization procedure always leads to a reduction in fuel consumption. The improvement that can be achieved via trim optimization increases when  $\delta$  is increased, although this tendency seems to stabilize for  $\delta > 5\%$ . According to the results of this model, improvements exceeding 2% in fuel consumption can be achieved by applying the model for trim optimization to the selected vessel. It should be noted that trim optimization can be performed at near to zero cost on board, since it does not require the installation of any additional equipment.

#### 7.4.3.4 Conclusions

Authors of [87] have shown how data driven models (BBMs) can outperform state-of-the-art numerical models (WBM) which exploits the physical knowledge of the system in the task of predicting the fuel consumption of a naval propulsion plant. Based on these models new approaches for modeling the system have been developed, namely the GBMs, which are able to exploit the advantages of two philosophies: GBMs are able to obtain the same performances of the BBMs but requiring less historical data thanks to the knowledge embedded in the WBM. The proposed methodologies have been tested on real world historical data collected from a real vessel during two years of on board sensors data acquisitions, and the physical plausibility of the models has

been checked through a feature ranking process. Feature ranking allowed improving the understanding of BBMs and GBMs as for these model physical principles are only partly accounted for. Thanks to the high accuracy and physical plausibility of the developed models, the authors have been able to propose a trim optimization technique which exploits the predictive power of the proposed models for the online selection of the best configuration of the trim for reducing the fuel consumption. Results have shown to be very promising and they should be further verified by implementing the proposal on the onboard system of a vessel.

## 7.5. Future of data science and advanced analytics

According to Academia [89–92] and Industry [93–96] there are many Data Science and Advanced Analytics technologies that should be better studied and engineered because of their potential benefits for the shipping energy systems domain.

In the last years Data Science and Advanced Analytics are experiencing a fast process of commodification [97–102]. This characterization is on the interest of big IT companies, but it correctly reflects the current industrialization of Data Science and Advanced Analytics also in the field of shipping energy systems. This phenomenon means that these systems and products are reaching the society at large and, therefore, the trustworthiness related to the use of these tools cannot be ignored any longer. Designing technologies from this human-centered perspective means incorporating human-relevant requirements such as safety, fairness, privacy, and interpretability, but also considering broad societal issues such as ethics and legislation. These are essential aspects to foster the acceptance of Data Science and Advanced Analytics technologies in a human oriented environment and the shipping one, as well as to be able to comply with an evolving legislation concerning the impact of digital technologies on ethically and privacy sensitive matters. These technologies were not originally conceived with an eye on ethical issues but they were simply trying to emulate certain aspects of biological intelligence. It might also be argued that one of the aspects of the outcome of human biological intelligence is precisely unethical behavior. It is true, though, that ethics do only come into play in social interaction, where different human intelligence communicate and interact with each other. Right now, we find ourselves at a crossroad in the development intelligent entities that are beginning to become tightly interwoven to the shipping environment, in the form, for example, of intelligent assistants for maintenance. Unsurprisingly, their societal impact is coming to the fore of public discussion. For this purpose, these technologies are now requested to satisfy some additional requirements such as Privacy, Fairness, Safety, Security, Reliability, Interpretability, and Explainability. The problem of learning from data while preserving the privacy of individual observations has a long history and spans over multiple disciplines. One way to preserve privacy is to corrupt the learning procedure with noise without destroying the

information that we want to extract. Another way is to exploit the data in a federated way, leaving the data in the hand of the data generator (on the edge) centralizing only an aggregated information. Safety, Security, and Reliability are the property of Data Science and Advanced Analytics to be able to provide robust answers and suggestions. Recent deep learning algorithms (e.g., computer vision) have even surpassed human performances on some well-defined benchmark datasets. It has thus been extremely surprising to discover that such algorithms can be easily fooled by adversarial examples, that are, imperceptible, adversarial perturbations that mislead these systems into perceiving things that are not there. This undermined the safety and security properties of such algorithms and a large number of stakeholders have shown interest in understanding the risks associated to their misuses, to develop proper mitigation strategies and incorporate them in their product. Finally, we have to highlight that one of the legal bottlenecks hampering the application of Data Science and Advanced Analytics to real problems in the social domain is the “right to explanation” granted to citizens. Such requirement is in direct course of collision with the limitations of many Data Science and Advanced Analytics technologies in terms of interpretability and explainability. These issues have late come to the forefront of researchers, mostly due to the widespread development and application of Deep Learning methods in systems with societal impact. As they amplify shallow neural networks, it comes as no surprise that Deep Learning may become an extreme case of black box model, further reducing their interpretability and explainability.

Another problem of Data Science and Advanced Analytics technologies, which limits their adoption, is the difficulty of choosing the right tool for a specific problem [54,103–109]. The series of no-free-lunch theorems [110–114] ensure that there is not, and there will never be, a single tool able to efficiently solve all tasks and for this reason it is required to develop meta-tools able to select the right one for the specific case under exam. For example, in ML there is the need to automatically select the best learning algorithm or the best hyperparameter for a specific algorithm. In optimization we need to select the best language for describing the problem, or the best solver, or the best optimization strategy, or the best parameters of the optimizer. For this reason researchers are starting to combine different tools for building the so called meta-algorithms. In ML optimization strategies of algorithms are exploited not for solving the specific problem under exam but for selection of the best tool: for example, it is constructed an optimized or a predictive model to guide the selection of a particular tool that should be better suited for the specific problem under exam. In optimization ML models are exploited to determine the best solver to exploit, the best optimization strategy, the best description language, of to guide the optimizer toward better local minima. The final focus of all these researches is to involve the human in the loop for solving a specific task as less as possible so to make the tools accessible in as many contexts as possible with minimal knowledge and intervention.

As also shown in Section 7.3, while Descriptive, Diagnostic, and Predictive analytics are nowadays quite exploited both in research and in practice, less examples of Prescriptive Analytics can be found. In fact Prescriptive Analytics is the effort to fully automatize the process of taking decisions and actions starting from the data about the problem with no human intervention making specific processes (e.g., maintenance or fuel optimization) autonomous [26,27,115–118]. On the one side, this process is limited by the specific domain of the shipping energy system, which requires (because of the legislation or because of the contracts) that the final decision should be undertaken by a human operator which takes responsibility for that choice (and this is why Visual Analytics is so important). However, on the other side, there is a technological limitation: Prescriptive Analytics requires the knowledge of multiple aspects of artificial intelligence and the presence of multiple data sources which are not always available. For example, to model constraints and preferences of the operators, we need to exploit data in the form of ontology describing the context (and in the shipping energy systems there is still a large gap in this sense) and we need to exploit data and information which is not structured (we need to use optical recognition or audio recognition tool to extract the information from the human operator reports and we need to process it with natural language processing tools) to achieve practical results. This process requires a big effort in research to adapt and improve the current tools to the ship energy systems context but also a big investment from the companies in developing internally the skills required to adopt these tools. While for simpler analytics this process started already many years ago, for more advanced analytics this process is still in its early stages and more effort is required to fill the current gaps.

## References

- [1] M. Stopford, *Maritime Economics 3e*, Routledge, 2008.
- [2] P. Skerry, *Counting on the Census?: Race, Group Identity, and the Evasion of Politics*, vol. 56, Brookings Institution Press, 2000.
- [3] M. Siegler, Eric Schmidt: every 2 days we create as much information as we did up to 2003, Techcrunch, Beschikbaar via: <http://techcrunch.com/2010/08/04/schmidt-data>, 2010.
- [4] James Fisher and Sons plc, Mimic condition monitoring software, <http://www.jfmimic.co.uk>, 2021. (Accessed 1 April 2021).
- [5] J. R. Company, Japan radio co., ltd., <http://www.jrc.co.jp/eng/>, 2021. (Accessed 1 April 2021).
- [6] R. Perez, T. Cebecauer, M. Šúri, Semi-empirical satellite models, *Solar Energy Forecasting and Resource Assessment (2013)* 21–48.
- [7] N. R. Laboratory, Naval research laboratory, <https://www.nrl.navy.mil/>, 2021. (Accessed 1 April 2021).
- [8] N. W. Service, National weather service, <https://www.weather.gov/>, 2021. (Accessed 1 April 2021).
- [9] M. Office, Met Office, <https://www.metoffice.gov.uk/>, 2021. (Accessed 1 April 2021).
- [10] F. B.V., Offshore weather forecasting services, <https://www.fugro.com/our-services/marine-asset-integrity/monitoring-and-forecasting/offshore-weather-forecasting-services>, 2020. (Accessed 1 July 2020).
- [11] S. Ltd., Spire weather solutions, <https://spire.com/weather/>, 2020. (Accessed 1 July 2020).
- [12] E. S. Agency, Copernicus open access hub, <https://scihub.copernicus.eu/>, 2020. (Accessed 1 July 2020).

- [13] P.A. Janssen, S. Abdalla, H. Hersbach, J.R. Bidlot, Error estimation of buoy, satellite, and model wave height data, *Journal of Atmospheric and Oceanic Technology* 24 (2007) 1665–1677.
- [14] J.M. Giron-Sierra, J.F. Jimenez, State-of-the-Art of Wave Measurement for Ship Motion Prediction, vol. 43, IFAC, 2010.
- [15] P. Kasinatha Pandian, O. Emmanuel, J.P. Ruscoe, J.C. Side, R.E. Harris, S.A. Kerr, C.R. Bullen, An overview of recent technologies on wave and current measurement in coastal and marine applications, *Journal of Oceanography and Marine Science* 1 (2010) 1–10.
- [16] G. Ludeno, F. Raffa, F. Soldovieri, F. Serafino, X-band radar for the monitoring of sea waves and currents: a comparison between medium and short radar pulses, *Geoscientific Instrumentation, Methods and Data Systems Discussions* (2017) 1–11.
- [17] G. Ludeno, F. Raffa, F. Soldovieri, F. Serafino, Proof of feasibility of the sea state monitoring from data collected in medium pulse mode by a X-band wave radar system, *Remote Sensing* 10 (2018).
- [18] A. Joseph, Chapter 13 – conclusions, in: A. Joseph (Ed.), *Measuring Ocean Currents*, Elsevier, 2014, pp. 397–417.
- [19] P.S. Bell, J.C. Osler, Mapping bathymetry using X-band marine radar data recorded from a moving vessel, *Ocean Dynamics* 61 (2011) 2141–2156.
- [20] V. Dhar, Data science and prediction, *Communications of the ACM* 56 (2013) 64–73.
- [21] W. Van Der Aalst, Data science in action, in: *Process Mining*, 2016.
- [22] F. Provost, T. Fawcett, Data science and its relationship to big data and data-driven decision making, *Big Data* 1 (2013) 51–59.
- [23] F. Provost, T. Fawcett, *Data Science for Business: What You Need to Know About Data Mining and Data-Analytic Thinking*, O’Reilly Media, Inc., 2013.
- [24] N.A. Waller, S.E. Fawcett, Data science, predictive analytics, and big data: a revolution that will transform supply chain design and management, *Journal of Business Logistics* 34 (2013) 77–84.
- [25] N. Elgendy, A. Elragal, Big data analytics: a literature review paper, in: *Industrial Conference on Data Mining*, 2014.
- [26] R. Soltanpoor, T. Sellis, Prescriptive analytics for big data, in: *Australasian Database Conference*, 2016.
- [27] D. Bertsimas, N. Kallus, From predictive to prescriptive analytics, arXiv preprint, arXiv:1402.5481, 2014.
- [28] T. Maydon, The 4 types of data analytics, <https://www.kdnuggets.com/2017/07/4-types-data-analytics.html>, 2017. (Accessed 1 July 2020).
- [29] Ø.J. Rødseth, L.P. Perera, B. Mo, Big data in shipping—challenges and opportunities, [https://sintef.brage.unit.no/sintef-xmlui/bitstream/handle/11250/2390646/34\\_Rodsoth.pdf](https://sintef.brage.unit.no/sintef-xmlui/bitstream/handle/11250/2390646/34_Rodsoth.pdf), 2020. (Accessed 1 July 2020).
- [30] P. Russom, Big data analytics, TDWI best practices report, fourth quarter 19 (2011) 1–34.
- [31] K. Kambatla, G. Kollias, V. Kumar, A. Grama, Trends in big data analytics, *Journal of Parallel and Distributed Computing* 74 (2014) 2561–2573.
- [32] S. LaValle, E. Lesser, R. Shockley, M.S. Hopkins, N. Kruschwitz, Big data, analytics and the path from insights to value, *MIT Sloan Management Review* 52 (2011) 21–32.
- [33] C.W. Tsai, C.F. Lai, H.C. Chao, A.V. Vasilakos, Big data analytics: a survey, *Journal of Big Data* 2 (2015) 1–32.
- [34] I. Zaman, K. Pazouki, R. Norman, S. Younessi, S. Coleman, Challenges and opportunities of big data analytics for upcoming regulations and future transformation of the shipping industry, *Procedia Engineering* 194 (2017) 537–544.
- [35] A. Coraddu, L. Oneto, F. Baldi, D. Anguita, Vessels fuel consumption: a data analytics perspective to sustainability, in: *Soft Computing for Sustainability Science*, 2018.
- [36] S. Kandel, J. Heer, C. Plaisant, J. Kennedy, F. Van Ham, N.H. Riche, C. Weaver, B. Lee, D. Brodbeck, P. Buono, Research directions in data wrangling: visualizations and transformations for usable and credible data, *Information Visualization* 10 (2011) 271–288.
- [37] I.G. Terrizzano, P.M. Schwarz, M. Roth, J.E. Colino, Data wrangling: the challenging journey from the wild to the lake, in: *Biennial Conference on Innovative Data Systems Research*, 2015.
- [38] T. Furche, G. Gottlob, L. Libkin, G. Orsi, N.W. Paton, Data wrangling for big data: challenges and opportunities, in: *International Conference on Extending Database Technology*, 2016.

- [39] F. Endel, H. Piringer, Data wrangling: making data useful again, *IFAC-PapersOnLine* 48 (2015) 111–112.
- [40] T. Rattenbury, J.M. Hellerstein, J. Heer, S. Kandel, C. Carreras, *Principles of Data Wrangling: Practical Techniques for Data Preparation*, O'Reilly Media, Inc., 2017.
- [41] C.C. Aggarwal, *Data Mining: the Textbook*, Springer, 2015.
- [42] A. Azzalini, B. Scarpa, *Data Analysis and Data Mining: An Introduction*, OUP, USA, 2012.
- [43] A. Rajaraman, J.D. Ullman, *Mining of Massive Datasets*, Cambridge University Press, 2011.
- [44] M. Kantardzic, *Data Mining: Concepts, Models, Methods, and Algorithms*, John Wiley & Sons, 2011.
- [45] J. Wang, *Data Mining: Opportunities and Challenges*, Idea Group Pub., 2003.
- [46] S. Shalev-Shwartz, S. Ben-David, *Understanding Machine Learning: From Theory to Algorithms*, Cambridge University Press, 2014.
- [47] M. Mohri, A. Rostamizadeh, A. Talwalkar, *Foundations of Machine Learning*, MIT Press, 2018.
- [48] J. Alzubi, A. Nayyar, A. Kumar, Machine learning from theory to algorithms: an overview, *Journal of Physics Conference Series* 1142 (2018) 012012.
- [49] E. Alpaydin, *Introduction to Machine Learning*, MIT Press, 2020.
- [50] C.M. Bishop, *Pattern Recognition and Machine Learning*, Springer, 2006.
- [51] I. Goodfellow, Y. Bengio, A. Courville, *Deep Learning*, MIT Press, 2016.
- [52] J. Schmidhuber, Deep learning in neural networks: an overview, *Neural Networks* 61 (2015) 85–117.
- [53] R.H. Shumway, D.S. Stoffer, *Time series analysis and its applications*, *Studies in Informatics and Control* 9 (2000) 375–376.
- [54] L. Oneto, *Model Selection and Error Estimation in a Nutshell*, Springer, 2020.
- [55] N.J. Nilsson, *Principles of Artificial Intelligence*, Morgan Kaufmann, 2014.
- [56] A. Barr, E.A. Feigenbaum, *The Handbook of Artificial Intelligence*, Volume 2, Butterworth–Heinemann, 2014.
- [57] S. Russell, P. Norvig, *Artificial Intelligence: a Modern Approach*, Pearson, 2002.
- [58] V. Lifschitz, *Answer Set Programming*, Springer, 2019.
- [59] P. Haslum, N. Lipovetzky, D. Magazzeni, C. Muise, An introduction to the planning domain definition language, *Synthesis Lectures on Artificial Intelligence and Machine Learning* 13 (2019) 1–187.
- [60] J. Preece, Y. Rogers, H. Sharp, D. Benyon, S. Holland, T. Carey, *Human–Computer Interaction*, Addison–Wesley Longman Ltd., 1994.
- [61] J. Preece, H. Sharp, Y. Rogers, *Interaction Design: Beyond Human–Computer Interaction*, John Wiley & Sons, 2015.
- [62] A. Dix, A.J. Dix, J. Finlay, G.D. Abowd, R. Beale, *Human–Computer Interaction*, Pearson Education, 2003.
- [63] S.K. Card, *The Psychology of Human–Computer Interaction*, CRC Press, 2018.
- [64] M.G. Helander, *Handbook of Human–Computer Interaction*, Elsevier, 2014.
- [65] R. Arias-Hernández, J. Dill, B. Fisher, T.M. Green, Visual analytics and human–computer interaction, *Interactions* 18 (2011) 51–55.
- [66] R. Lu, O. Turan, E. Boulougouris, C. Banks, A. Incecik, A semi-empirical ship operational performance prediction model for voyage optimization towards energy efficient shipping, *Ocean Engineering* 110 (2015) 18–28.
- [67] K. Wang, X. Yan, Y. Yuan, X. Jiang, X. Lin, R. Negenborn, Dynamic optimization of ship energy efficiency considering time-varying environmental factors, *Transportation Research Part D, Transport and Environment* 62 (2018) 685–698.
- [68] K. Wang, X. Yan, Y. Yuan, F. Li, Real-time optimization of ship energy efficiency based on the prediction technology of working condition, *Transportation Research Part D, Transport and Environment* 46 (2016) 81–93.
- [69] L.T. Le, G. Lee, K.-S. Park, H. Kim, Neural network-based fuel consumption estimation for container ships in Korea, *Maritime Policy and Management* 47 (2020) 615–632.
- [70] O. Soner, E. Akyuz, M. Celik, Use of tree based methods in ship performance monitoring under operating conditions, *Ocean Engineering* 166 (2018) 302–310.
- [71] A. Pagoropoulos, A. Moller, T. McAloone, Applying multi-class support vector machines for performance assessment of shipping operations: the case of tanker vessels, *Ocean Engineering* 140 (2017) 1–6.

- [72] A. Parkes, A. Sobey, D. Hudson, Physics-based shaft power prediction for large merchant ships using neural networks, *Ocean Engineering* 166 (2018) 92–104.
- [73] X. Xu, Z. Zhao, X. Xu, J. Yang, L. Chang, X. Yan, G. Wang, Machine learning-based wear fault diagnosis for marine diesel engine by fusing multiple data-driven models, *Knowledge-Based Systems* 190 (2020) 105324.
- [74] R. Yan, S. Wang, Y. 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 (2020) 101930.
- [75] M. Jeon, Y. Noh, Y. Shin, O. Lim, I. Lee, D. Cho, Prediction of ship fuel consumption by using an artificial neural network, *Journal of Mechanical Science and Technology* 32 (2018) 5785–5796.
- [76] X. Yan, K. Wang, Y. Yuan, X. Jiang, R.R. Negenborn, Energy-efficient shipping: an application of big data analysis for optimizing engine speed of inland ships considering multiple environmental factors, *Ocean Engineering* 169 (2018) 457–468.
- [77] C. Zhang, D. Zhang, M. Zhang, W. Mao, Data-driven ship energy efficiency analysis and optimization model for route planning in ice-covered Arctic waters, *Ocean Engineering* 186 (2019) 106071.
- [78] S. Wang, B. Ji, J. Zhao, W. Liu, T. Xu, Predicting ship fuel consumption based on lasso regression, *Transportation Research Part D, Transport and Environment* 65 (2018) 817–824.
- [79] Y.B. Farag, A.I. Ölçer, The development of a ship performance model in varying operating conditions based on ANN and regression techniques, *Ocean Engineering* 198 (2020).
- [80] B. Yoo, J. Kim, Probabilistic modeling of ship powering performance using full-scale operational data, *Applied Ocean Research* 82 (2019) 1–9.
- [81] A. Coraddu, L. Oneto, D. Iardi, S. Stoumpos, T. Gerasimos, Marine dual fuel engines monitoring in the wild through weakly supervised data analytics, *Engineering Applications of Artificial Intelligence* (2021).
- [82] S. Stoumpos, G. Theotokatos, E. Boulougouris, D. Vassalos, I. Lazakis, G. Livanos, Marine dual fuel engine modelling and parametric investigation of engine settings effect on performance–emissions trade-offs, *Ocean Engineering* 157 (2018) 376–386.
- [83] A. Coraddu, L. Oneto, F. Baldi, F. Cipollini, M. Atlar, S. Savio, Data-driven ship digital twin for estimating the speed loss caused by the marine fouling, *Ocean Engineering* 186 (2019) 106063.
- [84] I.T. Committee, Ships and Marine Technology Measurement of Changes in Hull and Propeller Performance – Part 2: Default Method, Standard, International Organization for Standardization, Geneva, CH, 2016.
- [85] E. Glaessgen, D. Stargel, The digital twin paradigm for future NASA and US air force vehicles, in: 53rd Structures, Structural Dynamics and Materials Conference, 2012.
- [86] S. Boschert, R. Rosen, Digital twin – the simulation aspect, in: *Mechatronic Futures*, 2016.
- [87] A. Coraddu, L. Oneto, F. Baldi, D. Anguita, Vessels fuel consumption forecast and trim optimisation: a data analytics perspective, *Ocean Engineering* 130 (2017) 351–370.
- [88] A. Coraddu, M. Figari, S. Savio, D. Villa, A. Orlandi, Integration of seakeeping and powering computational techniques with meteo-marine forecasting data for in-service ship energy assessment, in: *Developments in Maritime Transportation and Exploitation of Sea Resources*, 2013.
- [89] D.F. Dominković, I. Bačeković, A.S. Pedersen, G. Krajačić, The future of transportation in sustainable energy systems: opportunities and barriers in a clean energy transition, *Renewable and Sustainable Energy Reviews* 82 (2018) 1823–1838.
- [90] S.J. Davis, N.S. Lewis, M. Shaner, S. Aggarwal, D. Arent, I.L. Azevedo, S.M. Benson, T. Bradley, J. Brouwer, Y.M. Chiang, Net-zero emissions energy systems, *Science* 360 (2018).
- [91] G. Zou, Intelligent design and operation of ship energy systems combining big data and AI, in: *Conference on Computer and IT Applications in the Maritime Industries*, 2018.
- [92] F. Ahlgren, *Reducing Ships' Fuel Consumption and Emissions by Learning from Data*, Linnaeus University Press, 2018.
- [93] Gartner, Top trends on the Gartner hype cycle for artificial intelligence, <https://www.gartner.com/smarterwithgartner/top-trends-on-the-gartner-hype-cycle-for-artificial-intelligence-2019/>, 2019. (Accessed 1 July 2020).



- [94] Gartner, 5 trends appear on the Gartner hype cycle for emerging technologies, <https://www.gartner.com/smarterwithgartner/5-trends-appear-on-the-gartner-hype-cycle-for-emerging-technologies-2019/>, 2019. (Accessed 1 July 2020).
- [95] A. Love, AI in shipping: areas to watch in 2020, <https://www.ship-technology.com/features/ai-in-shipping/>, 2020. (Accessed 1 July 2020).
- [96] DigitalShip, Artificial intelligence applied to vessel power systems, <https://thedigitalship.com/news/maritime-software/item/5272-artificial-intelligence-applied-to-vessel-power-systems>, 2020. (Accessed 1 July 2020).
- [97] D. Bacciu, B. Biggio, P.J.G. Lisboa, J.D. Martin Guerrero, L. Oneto, A. Vellido Alcacena, Societal issues in machine learning: when learning from data is not enough, in: European Symposium on Artificial Neural Networks, Computational Intelligence and Machine Learning, 2019.
- [98] L. Floridi, Establishing the rules for building trustworthy AI, *Nature Machine Intelligence* 1 (2019) 261–262.
- [99] M. Brundage, S. Avin, J. Wang, H. Belfield, G. Krueger, G. Hadfield, H. Khlaaf, J. Yang, H. Toner, R. Fong, Toward trustworthy AI development: mechanisms for supporting verifiable claims, arXiv preprint, arXiv:2004.07213, 2020.
- [100] D. Danks, The value of trustworthy AI, in: AAAI/ACM Conference on AI, Ethics and Society, 2019.
- [101] A. Kumar, T. Braud, S. Tarkoma, P. Hui, Trustworthy AI in the age of pervasive computing and big data, arXiv preprint, arXiv:2002.05657, 2020.
- [102] C.S. Wickramasinghe, D.L. Marino, J. Grandio, M. Manic, Trustworthy AI development guidelines for human system interaction, in: International Conference on Human System Interaction, 2020.
- [103] X. He, K. Zhao, X. Chu, Automl: a survey of the state-of-the-art, arXiv preprint, arXiv:1908.00709, 2019.
- [104] I. Guyon, L. Sun-Hosoya, M. Boullé, H.J. Escalante, S. Escalera, Z. Liu, D. Jajetic, B. Ray, M. Saeed, M. Sebag, Analysis of the automl challenge series, in: Automated Machine Learning, 2019.
- [105] J. Liang, E. Meyerson, B. Hodjat, D. Fink, K. Mutch, R. Miikkulainen, Evolutionary neural automl for deep learning, in: Genetic and Evolutionary Computation Conference, 2019.
- [106] A. Truong, A. Walters, J. Goodsitt, K. Hines, C.B. Bruss, R. Farivar, Towards automated machine learning: evaluation and comparison of automl approaches and tools, in: International Conference on Tools with Artificial Intelligence, 2019.
- [107] M. Feurer, F. Hutter, Towards further automation in automl, in: International Conference on Machine Learning, 2018.
- [108] D. Wang, P. Ram, D.K.I. Weidele, S. Liu, M. Muller, J.D. Weisz, A. Valente, A. Chaudhary, D. Torres, H. Samulowitz, AutoAI: automating the end-to-end AI lifecycle with humans-in-the-loop, in: International Conference on Intelligent User Interfaces Companion, 2020.
- [109] H.H. Hoos, Automated artificial intelligence (AutoAI), ADA Research Group Technical Report TR-2018, 2018.
- [110] D.H. Wolpert, The lack of a priori distinctions between learning algorithms, *Neural Computation* 8 (1996) 1341–1390.
- [111] D.H. Wolpert, W.G. Macready, No free lunch theorems for optimization, *IEEE Transactions on Evolutionary Computation* 1 (1997) 67–82.
- [112] D.H. Wolpert, W.G. Macready, No free lunch theorems for search, Technical Report SFI-TR-95-02-010, Santa Fe Institute, 1995.
- [113] D.H. Wolpert, The supervised learning no-free-lunch theorems, in: *Soft Computing and Industry*, 2002.
- [114] S.P. Adam, S.A.N. Alexandropoulos, P.M. Pardalos, M.N. Vrahatis, No free lunch theorem: a review, in: *Approximation and Optimization*, 2019.
- [115] K. Lepenioti, A. Bousdekis, D. Apostolou, G. Mentzas, Prescriptive analytics: literature review and research challenges, *International Journal of Information Management* 50 (2020) 57–70.
- [116] C. Gröger, H. Schwarz, B. Mitschang, Prescriptive analytics for recommendation-based business process optimization, in: *International Conference on Business Information Systems*, 2014.
- [117] D. Bertsimas, B. Van Parys, Bootstrap robust prescriptive analytics, arXiv preprint, arXiv:1711.09974, 2017.
- [118] L. Berk, D. Bertsimas, A.M. Weinstein, J. Yan, Prescriptive analytics for human resource planning in the professional services industry, *European Journal of Operational Research* 272 (2019) 636–641.