

SPAMDA: Software for Pre-processing and Analysis of Meteorological DATA to build datasets

Antonio M. Gómez-Orellana^a, Juan C. Fernández^{a,*}, Manuel Dorado-Moreno^a,
Pedro A. Gutiérrez^a, César Hervás-Martínez^a

^a*Department of Computer Science and Numerical Analysis, University of Cordoba,
Córdoba, Spain.*

Abstract

Meteorological data play an important role for the comprehension of the environment, and they are extensively used to perform environmental learning. Machine Learning (ML) techniques have become a valuable support in many research areas, improving the performance of traditional statistical procedures. However, such techniques require datasets containing information related to the topic under study, which are not always available in an appropriate format. Preparing these datasets implies a lot of time and effort by the researchers. This paper presents a software tool for creating datasets using meteorological observations from two well-known sources of information: the *National Data Buoy Center* (NDBC) and the *National Centers for Environmental Prediction* (NCEP)/*National Center for Atmospheric Research* (NCAR) *Reanalysis Project*. The datasets created by the software are ready to be used as inputs for ML techniques in prediction tasks (classification or regression). As the designed software is able to simplify the creation of the datasets and reduce the time needed for this task, it prevents researchers from performing repetitive technical work, allowing them to concentrate on the study of the meteorological aspects of the data. Therefore, researchers can benefit from it in order to achieve a more efficient exploitation and protection of the environment.

*Corresponding author.

Email addresses: `i32goora@uco.es` (Antonio M. Gómez-Orellana), `jfcaballero@uco.es` (Juan C. Fernández), `manuel.dorado@uco.es` (Manuel Dorado-Moreno), `pagutierrez@uco.es` (Pedro A. Gutiérrez), `chervas@uco.es` (César Hervás-Martínez)

Keywords: Meteorological data, Reanalysis data, Pre-processing, Creating datasets, Prediction tasks, Marine energy

1. Introduction

A better understanding of the environment is of vital importance for science, contributing not only to more efficient exploitation of natural resources but also to the development of new strategies aimed at its protection. In that sense, meteorological observations provide an essential and valuable source of information which is widely used by researchers to address environmental learning, comprehension, prediction and conservation in numerous oceanic and atmospheric studies of a wide variety of areas (e.g. climate change, agriculture, energy, etc.). Some specific examples of the diversity of fields in which meteorological data can be used in are, among others: directional analysis of sea storms [1], wind power ramp events prediction [2], offshore wind energy assessment [3], study of the responses exhibited by plankton to fluid motions [4], trends in solar radiation [5] or simulation of extreme near shore sea conditions [6]. All these studies require a prior data collection and its adaptation to a specific format that allows the interpretation of them.

On the other hand, special purpose software is usually developed to help researchers to advance in their studies related to energy and environmental modelling, becoming a great support for decision-making in the exploitation and protection of the environment. In [7], a software tool for designing solar water heating systems is developed, which simulates different situations and finds the best technical and economical solution. In [8], an integrated simulation tool for the optimum design of bifacial solar panel with reflectors is presented. This tool can also be used to analyse the performance of the solar cells. A framework for data integration of offshore wind farms is implemented in [9] in order to facilitate data exchange and improve operation and maintenance practices. In [10], a tool based on *Fatigue, Aerodynamics, Structures, and Turbulence* (FAST) code for performing design optimisation of offshore wind turbines is presented, which can

analyse a massive group of cases by means of its parametric design capability. A design tool for architects based on *Lighting and Thermal* (LT) method is used in [11], with the purpose of comparing and optimising design solutions in terms of energy and comfort performance. Raabe et al. [12] developed two software tools, *Model of Equilibrium of Bay Beaches* (MEPBAY) and *Coastal Modeling System* (SMC), to support different operational levels of headland-bay beach in coastal engineering projects, and a software tool for the creation of a *Typical Meteorological Year* (TMY) that can be used for designing of solar energy systems is presented in [13].

Marine energy prediction is currently a hot topic where meteorological data is used in. Marine Renewable Energy (MRE) is one of the most important renewable and sustainable energy sources available in our environment, and it includes ocean thermal energy, marine tidal current energy and wave energy, among others. Its benefits and great potential [14] make it one of the most relevant natural resources, playing a crucial role not only in the reduction of the emission of greenhouse gases but also in all other aspects involved in the difficult challenge of the transition to a low carbon footprint society [15]. Wave energy exhibits a more stable power supply than wind energy and even solar energy. In recent years *Wave Energy Converters* (WECs) [16] have been developed to transform this wave energy into electricity. WECs are mechanical devices that convert kinetic energy into electrical energy by means of either the vertical oscillation of waves or the linear motion of them. Nevertheless, waves are difficult to be characterised due to their stochastic nature, because of the influence of a large number of environmental factors that exert on them [17]. As a consequence of this complexity, many aspects of WEC design, deployment and operation [18, 19, 20] need a proper prediction of waves, in order to maximise the wave energy extraction. For this purpose WECs use wave *flux of energy* (F_e) which can be calculated from the two most important wave parameters in this regard: *significant wave height* (H_s) and *wave energy period* (T_e). Additionally, wave predictions are also helpful for designing offshore structures [21], operational works in the sea [22], etc.

Currently, and as a support to traditional study procedures, Machine Learning (ML) techniques [23, 24] are being widely used in numerous research fields related to classification, regression and optimisation tasks [25], obtaining significant improvements in the performance of the results. ML methodologies can be used not only by experienced data and computer scientists but also by other researchers. For example, the well-known *Waikato Environment for Knowledge Analysis* (WEKA) [26] software tool provides researchers with a wide collection of ML algorithms. In this way, ML techniques have been already applied to tackle wave characterisation, accurately estimating H_s and T_e parameters [27, 28], given that robustness of ML methods can tackle the previously explained difficulties in wave energy prediction. The problem is that, in order to apply these methods, it is essential to obtain datasets with relevant information about the issue under study, used to infer knowledge. Usually, these datasets are not publicly available in a friendly format, and their generation is the first step needed.

The information to create these datasets can be obtained from meteorological observations, but such information may be available in an inappropriate format and even contain missing values or measurements. Consequently, it is usually required to perform pre-processing tasks for improving the quality of the data, such as the replacement of missing values, outlier detection or data normalisation, among others. Furthermore, if more than one source of information is used to achieve a better characterisation of the problem under study [29, 30], then a merging or matching process have to be carried out by the researchers to manually create the datasets with the needed information. Moreover, depending on the subject and the ML technique to be applied, or even if the researcher considers other factors in order to improve the results or have more in-depth conclusions, the datasets would have to be updated afterwards. In summary, many important details and different intermediate steps have to be considered when creating suitable datasets for ML techniques, resulting in an extremely tedious task.

The main purpose of this paper is to present a software tool able to prevent

90 researchers from performing this repetitive work and greatly simplify all the
 steps involved in the creation of datasets. The meteorological data considered for
 the tool come from two well-known sources of information: the *National Oceanic
 and Atmospheric Administration* (NOAA) *National Data Buoy Center* (NDBC)
 [31] and the *National Centers for Environmental Prediction* (NCEP)/*National*
 95 *Center for Atmospheric Research* (NCAR) *Reanalysis Project* (NNRP or R1)
 [32, 33]. The software tool presented in this work is named SPAMDA (Software
 for Pre-processing and Analysis of Meteorological DAta to build datasets). As
 SPAMDA performs all this data processing, it reduces the time involving these
 tasks and allows the researchers focus on the study of the meteorological aspects
 100 of the observations. The datasets obtained are ready to be used as input for
 ML techniques in prediction tasks (classification or regression), although the
 researchers can use them for other purposes. These datasets contain one or
 more meteorological variables as inputs and one variable as target (variable to
 be predicted). The format of the generated datasets will be *Attribute-Relation*
 105 *File Format* (ARFF) [34], which is the one used by WEKA. Besides, the datasets
 can also be generated in *Comma-Separated Values* (CSV) format, enabling the
 researchers to use others tools.

Up to now, and to the best of our knowledge, this is the first software
 tool addressing the problem previously discussed, combining meteorological data
 110 from NDBC and NNRP. The advantages that SPAMDA offers to the researchers
 will be detailed in Section 3, although some of them are briefly summarised
 below:

- The generation of datasets becomes a very easy and customizable task, by
 means of the selection of different input parameters.
- 115 • It makes the researcher focus on oceanic and atmospheric studies, without
 having to worry about mechanical tasks.
- It provides information about the quality and quantity of the data.
- It avoids possible researcher errors in the intermediate steps of the process

of creation of the datasets.

- 120 • It includes different pre-processing tasks, such as normalisation and missing data recovery.
- It facilitates data management and well-organised storage of the datasets.
- Its modular design allows the implementation of new functional modules for managing meteorological data from others sources for renewable energy
- 125 research.
- It includes an user-friendly GUI, facilitating and greatly simplifying data management, and it is integrated with the Explorer environment of WEKA.
- It is multi-platform, and it can be used on any computer with Java regardless of the operating system.

130 This paper is organised as follows: Section 2 describes the sources of information used by SPAMDA for creating datasets. Section 3 describes in detail the features of the software tool. Section 4 shows a case study describing the use of SPAMDA in a practical approach. Section 5 provides the final conclusions and future work.

135 2. Meteorological data sources

The data provided by the above-mentioned sources of information of SPAMDA is described below:

- NDBC is a part of the *National Weather Service* (NWS). NDBC designs, develops, operates, and maintains a network of data collecting buoys (sta-
- 140 tions). The mission of the network is to collect real-time marine meteorological and oceanographic observations, such as H_s , dominant wave period, or wind speed and direction, among others.

The buoys maintained by NDBC are deployed in the coastal and offshore waters around oceans and seas, and they are equipped with assorted sen-

sors which allow them to perform different measurements. The information collected by the buoys is available on the NDBC website [35], and it is divided into different groups. One of them corresponds to standard meteorological information of the historical data collected by each buoy, which can be downloaded as annual text files and whose format was adopted by NDBC since January 2007 [36]. These files contain hourly measurements per day from 00:50 to 23:50 UTC (Universal Time Coordinated) and from 23:50 31th December of the previous desired year to 22:50 31th December of the desired year. In Table 1, a comprehensive measurement description and the corresponding units are provided as a summary for the reader. A fragment of one of these files, which contains the measurements collected during year 2017 by the buoy identified as *Station 46001* in NDBC, is shown in Fig. 1. Each column corresponds to a meteorological variable or attribute, and each row or instance corresponds to the values of the measurements collected by the buoy for each attribute at a specific date and time.

#YY	MM	DD	hh	mm	WDIR	WSPD	GST	WVHT	DPD	APD	MWD	PRES	ATMP	WTMP	DEWP	VIS	TIDE
#yr	mo	dy	hr	mn	degT	m/s	m/s	m	sec	sec	degT	hPa	degC	degC	degC	mi	ft
2016	12	31	23	50	279	6.4	7.3	2.41	12.90	6.50	999	1041.3	5.6	6.4	999.0	99.0	99.00
2017	01	01	00	50	291	6.3	7.3	2.13	7.14	6.08	999	1041.1	5.5	6.4	999.0	99.0	99.00
2017	01	01	01	50	293	5.3	6.6	2.39	7.69	6.64	999	1041.2	5.5	6.4	999.0	99.0	99.00
.
.
.
.
2017	12	31	20	50	999	3.0	4.4	4.98	12.90	8.57	201	1000.6	4.8	4.9	999.0	99.0	99.00
2017	12	31	21	50	999	3.8	6.3	4.64	10.00	8.55	150	1000.1	4.8	4.9	999.0	99.0	99.00
2017	12	31	22	50	999	3.4	5.2	4.40	12.90	8.40	200	998.9	4.7	4.9	999.0	99.0	99.00

Figure 1: A fragment of an annual text file of the *Station 46001*.

Note that the data collected by the network of buoys may be incomplete due to diverse circumstances such as the weather conditions in which the buoys have to operate, failures or malfunctioning elements of the buoys, among others. Accordingly, it may be the situation that some of the measurements are completely missing (missing date or instance) or partially missing (some measurements not recorded), by a buoy or by a set of buoys, once in a while or over a period of time. It may be also possible that the

Table 1: Measurements descriptions and units of each meteorological variable or attribute collected by the buoys.

Attribute	Units	Description
WDIR	degT	Wind direction (the direction the wind is coming from in degrees clockwise from true North) during the same period used for WSPD.
WSPD	m/s	Wind speed (m/s) averaged over an eight-minute period for buoys and a two-minute period for land stations. Reported Hourly.
GST	m/s	Peak 5 or 8 second gust speed (m/s) measured during the eight-minute or two-minute period.
WVHT	m	Significant wave height (meters) is calculated as the average of the highest one-third of all of the wave heights during the 20-minute sampling period.
DPD	sec	Dominant wave period (seconds) is the period with the maximum wave energy.
APD	sec	Average wave period (seconds) of all waves during the 20-minute period.
MWD	degT	The direction from which the waves at the dominant period (DPD) are coming. The units are degrees from true North, increasing clockwise, with North as 0 (zero) degrees and East as 90 degrees.
PRES	hPa	Sea level pressure (hPa). For C-MAN sites and Great Lakes buoys, the recorded pressure is reduced to sea level using the method described in NWS Technical Procedures Bulletin 291 (11/14/80).
ATMP	degC	Air temperature (Celsius degrees).
WTMP	degC	Sea surface temperature (Celsius degrees). For buoys the depth is referenced to the hull's waterline. For fixed platforms it varies with tide, but is referenced to, or near Mean Lower Low Water (MLLW).
DEWP	degC	Dewpoint temperature taken at the same height as the air temperature measurement.
VIS	nmi	Station visibility (nautical miles). Note that buoy stations are limited to reports from 0 to 1.6 nmi.
TIDE	ft	The water level in feet above or below MLLW.

measurements have been recorded at a time different from the expected one. This aspect have to be taken into account when creating the datasets. This casuistry is explained in detail in Appendix A.

• NNRP provides three-dimensional global reanalysis of numerous meteorological variables (e.g. air temperature, components South-North and West-East of wind speed, relative humidity, pressure, etc.), which is available monthly, daily and every 6 hours at 00 Z (Zulu time), 06 Z, 12 Z and 18 Z from 1948 on a global $2.5^\circ \times 2.5^\circ$ grid. Weather observations are from different sources, such as ships, satellites and radar, among others. Reanalysis data is created assimilating such observations using the same climate model throughout the entire reanalysis period in order to reduce the effects of modelling changes on climate statistics. Such information has become a substantial support of the needs of the research community, even more in locations where instrumental (real time) data is not available. The reanalysis data is available in the NNRP website [37], which it is accessible through different sections. Such data can be fully (a global $2.5^\circ \times 2.5^\circ$ grid) or partially (only the desired reanalysis nodes or sub-grid) downloaded as *Network Common Data Form* (NetCDF) files [38], a special binary format for representing scientific data, which provides a description of the file contents and also includes the spatial and temporal properties of the data. Each reanalysis file contains the values of a meteorological variable estimated by a mathematical model for each reanalysis node. For a better understanding, in Fig. 2 an approximate representation of a sub-grid containing six reanalysis nodes around the geographical location of a buoy (obtained from NDBC) is shown.

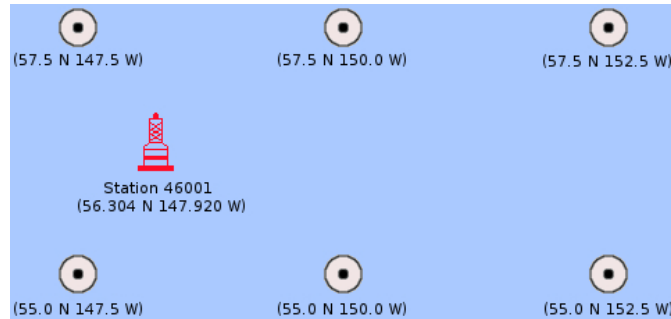


Figure 2: Example of a six sub-grid reanalysis nodes around the *Station 46001*.

Therefore, with both sources of information, which complement each other, and carrying out a matching process, SPAMDA will create datasets for prediction tasks. In this way, the dataset input variables will be one or more reanalysis variables from NNRP and one or more measurements from NDBC. The dataset output variable will always be one measurement from NDBC.

3. SPAMDA

SPAMDA combines meteorological information from NDBC and NNRP to obtain new datasets for oceanic and atmospheric studies. In order to do so, SPAMDA manages three different types of datasets, which will be described in detail in the following sections, but are briefly introduced bellow for giving the reader a better general understanding:

- *Intermediate datasets*: They contain the meteorological observations from NDBC.
- *Pre-processed datasets*: They are obtained as a result of pre-processing tasks performed on the intermediate datasets.
- *Final datasets*: Created by merging an intermediate or pre-processed dataset (which contain the information from NDBC) with the reanalysis data from NNRP. This procedure is referenced in SPAMDA as matching process and will be carried out according to the study to be performed (classification or regression).

SPAMDA consists of three main functional modules, whose main features, represented in Fig. 3, are the following:

- *Manage buoys data*: The aim of this module is to provide features for the management and analysis of the information related to the buoys from NDBC. This includes:
 1. Entering and updating the information of each buoy.
 2. Creation of intermediate datasets with the collected measurements.

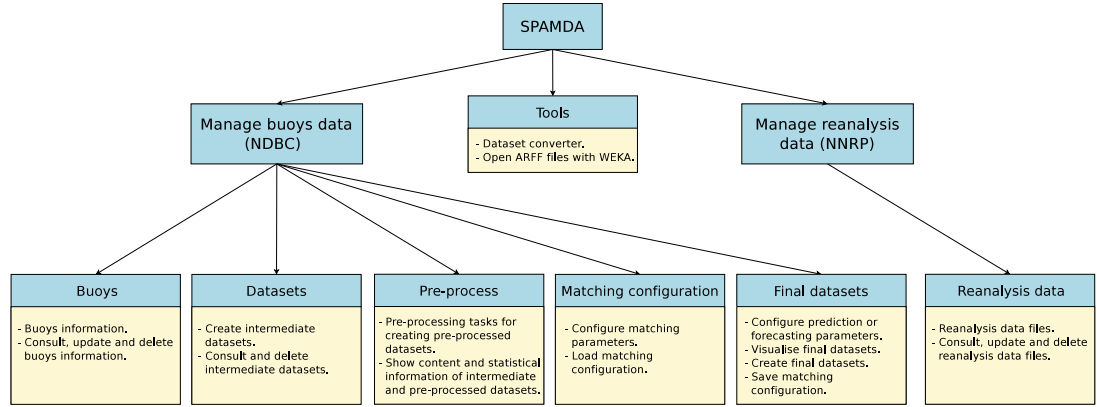


Figure 3: Brief outline of the functionality provided by SPAMDA.

- 220 3. Pre-processing tasks for obtaining the pre-processed datasets.
 4. Matching process to merge the information from NDBC and NNRP
 5. Creation of the final datasets accordingly to the ML technique to use
(classification or regression).
- 225 • *Manage reanalysis data*: This module is used for the management of the reanalysis data provided by the NNRP. In this way, the researchers can keep the reanalysis data files updated for their studies. Such files will be used, depending on the researchers needs, in the matching process when obtaining the final datasets.
 - 230 • *Tools*: This module includes features for converting intermediate or pre-processed datasets to ARFF or CSV format and for opening ARFF files with WEKA software.

In the following subsections each integrated functional module is described in detail.

3.1. Buoys

- 235 When a new buoy is included in SPAMDA the following information, which can be obtained from NDBC, is requested:

- *Station ID*: An alphanumeric identifier that allows easy identification of the buoy.
- *Description*: A short description of the buoy.
- 240 • *Latitude*: North or South geographical localisation (degrees) of the buoy.
- *Longitude*: West or East geographical localisation (degrees) of the buoy.
- *Measurements files*: The above-mentioned annual text files of the standard meteorological information collected by the buoy and downloaded from the NDBC website. This will be used for the creation of the intermediate
245 datasets. One file per year is expected.

For clarification, an example is presented in Fig. 4, where the buoy ID1 has three annual text files and the buoy ID2 has two annual text files.

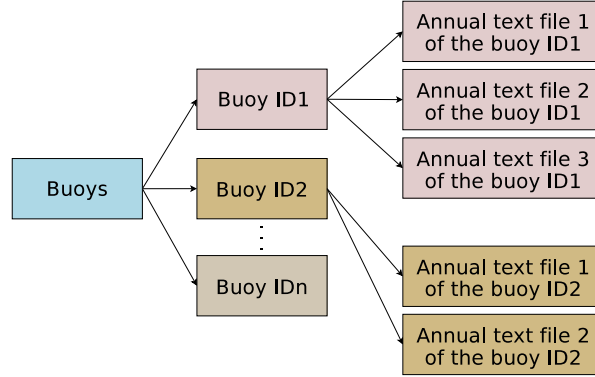


Figure 4: Example of entering two buoys with its annual text files.

3.2. Datasets

Once a buoy has been included as described in Section 3.1, it is possible
250 to create datasets with one or more annual text files, which are referenced in SPAMDA as intermediate datasets. In this module, the researchers can manage intermediate datasets of each buoy, which are the baseline for their studies, by creating new ones or deleting the unnecessary ones.

When an intermediate dataset is created, it is associated with its corresponding buoy. Besides, a summary of its content is also created, providing relevant information such as the number of instances, the dates of the first and last measurements, the annual text files included, and the missing and duplicated dates.

An example where three intermediate datasets have been created is presented in Fig. 5. The two intermediate datasets of the buoy ID1 contain meteorological data of different years, and the intermediate dataset of the buoy ID2 contains meteorological data of two years. For each buoy, as many intermediate datasets as needed can be created.

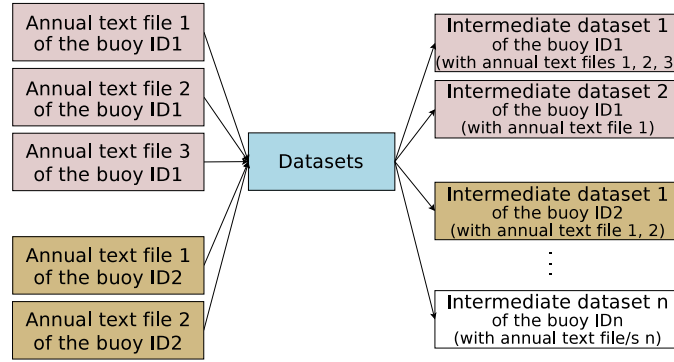


Figure 5: Example of the creation of the intermediate datasets.

3.3. Pre-process

Data pre-processing prepares the raw data (intermediate datasets) to be able to be treated correctly by ML algorithms. In this way, the quality of data can be improved prior to the learning phase, by applying pre-processing tasks (filters). The result will be referenced as pre-processed datasets.

SPAMDA provides several filters grouped in three categories, *Attribute*, *Instance* and *Recover missing data*, including the configuration of their parameters and a short description of them:

- *Attribute*: All these filters can be applied to the attributes (variables of the buoy from NDBC) of the intermediate dataset.

- 275
 – *Normalize* [39]: This filter normalises all numeric values of each attribute. The resulting values are by default in the interval $[0,1]$.
- *Remove* [40]: It removes an attribute or a range of them.
- *RemoveByName* [41]: It removes attributes based on a regular expression matched against their names.
- 280
 – *ReplaceMissingValues* [42]: For each attribute, all the missing values will be replaced by the average value of the attribute.
- *ReplaceMissingWithUserConstant* [43]: This filter replaces all the missing values of the attributes with an user-supplied constant value.
- *Instance*: All theses filters can be applied to the instances (hourly measurements of the buoy from NDBC) of the intermediate dataset.
- 285
 – *RemoveDuplicates* [44]: With this filter, all duplicated instances are removed.
- *RemoveWithValues* [45]: This filter removes all the instances that match the attribute and the value supplied by the user.
- 290
 – *SubsetByExpression* [46]: It removes all the instances which do not match a user-specified expression.
- *Recover missing data*: All these filters can be applied to the instances of the intermediate dataset.
- *Replace missing values with next nearest hour*: The missing values of each attribute are replaced with the next nearest non missing value.
- 295
 – *Replace missing values with previous nearest hour*: This filter replaces the missing values of each attribute with the previous nearest non missing value.
- *Replace missing values with next n hours mean*: The missing values of each attribute are replaced with the next n nearest non missing values mean, where n can be configured by the user.
- 300

- *Replace missing values with previous n hours mean*: This filter replaces the missing values of each attribute in the intermediate dataset with the previous n nearest non missing values mean.
- *Replace missing values with symmetric n hours mean*: The missing values of each attribute in the intermediate dataset are replaced with the n previous and n next non missing values mean.

305

SPAMDA allows the researchers to undo the last filter applied or to restore the initial content of the intermediate dataset. Besides, the content and relevant statistical information (number of instances with missing values, minimum and maximum values, mean and standard deviation) of the intermediate and the pre-processed datasets can be visualised in this module.

Fig. 6 shows an example where the intermediate datasets 1 and 2 of the buoy ID1 have been pre-processed, obtaining as a result the pre-processed dataset 1 of each one. The intermediate dataset 1 of the buoy ID2 has been also pre-processed. *Pre-processed dataset n* represents that the researchers can create as many pre-processed datasets as they consider opportune.

315

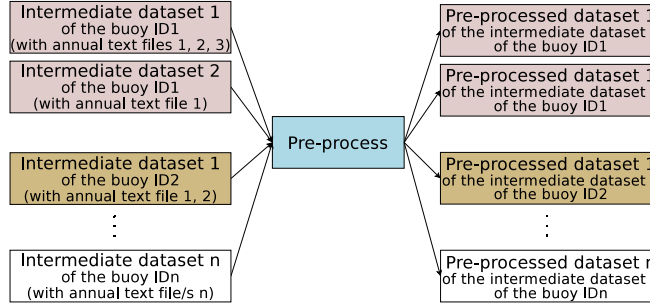


Figure 6: Example of the creation of pre-processed datasets.

3.4. Matching configuration

In order to merge and format the data provided by the two sources of information described in Section 2, it is necessary to carry out a matching process.

The matching procedure is performed using an intermediate or pre-processed

320

dataset, which includes the measurements collected by a buoy from NDBC, and the needed reanalysis data files from NNRP. Note that SPAMDA is able to manage the NetCDF binary format for handling the information stored in the reanalysis files.

325 Such process merges the information of both sources that match on time, but, given that the measurements of the buoys are hourly collected from 00:50 to 23:50 UTC, and the reanalysis data is available every 6 hours at 00 Z, 06 Z, 12 Z and 18 Z, the matching can only be carried every 6 hours (discarding the rest of measurements from the buoy data). Besides, and since there is still a difference
330 of 10 minutes, the matching with the reanalysis data will be performed with the nearest buoy measurement (before or after) within a maximum of 60 minutes of difference. Finally, the matched instances of both sources will form the final datasets.

Fig. 7 presents an example of matching with the measurements collected
335 during 2017 by *Station 46001* (NDBC) and the reanalysis data (NNRP) of the variable *pressure* for reanalysis nodes $57.5\text{ N} \times 147.5\text{ W}$ and $55.0\text{ N} \times 147.5\text{ W}$ in the same year. In this way, only the instances from both sources that are linked with arrows (highlighted in green colour) will be used in the creation of the final datasets. Although the reanalysis dates have been presented in a
340 human readable format, note that reanalysis dates are stored in hours from 01-01-1800, and they have to be transformed for comparison taking into account the time zone. Such transformation is automatically done by SPAMDA when matching the instances.

The reader can check in Appendix A 5 an example with a more complex
345 case of the procedure.

SPAMDA allows the researchers to perform a customisable matching process, for obtaining as many different versions of the same meteorological data as needed. Prediction tasks are based on the estimation of the output attribute using the information provided by the input attributes. Depending on the task,
350 the datasets must be prepared and configured differently:

#YY	#yr	mo	dy	hr	mn	WDIR	WSPD	GST	WVHT	DPD	APD	MWD	PRES	ATMP	WTMP	DEWP	VIS	TIDE
						degT	m/s	m/s	m	sec	sec	degT	hPa	degC	degC	degC	mi	ft
2016	12	31	23	50	279	6.4	7.3	2.41	12.90	6.50	999	1041.3	5.6	6.4	999.0	99.0	99.00	
2017	01	01	00	50	291	6.3	7.3	2.13	7.14	6.08	999	1041.1	5.5	6.4	999.0	99.0	99.00	
2017	01	01	01	50	293	5.3	6.6	2.39	7.69	6.64	999	1041.2	5.5	6.4	999.0	99.0	99.00	
2017	01	01	02	50	297	3.9	4.9	2.57	14.81	7.39	999	1041.1	5.5	6.4	999.0	99.0	99.00	
2017	01	01	03	50	298	5.7	6.7	1.83	7.69	5.80	999	1041.4	5.6	6.4	999.0	99.0	99.00	
2017	01	01	04	50	281	5.0	6.1	2.29	13.79	7.14	999	1041.5	5.6	6.4	999.0	99.0	99.00	
2017	01	01	05	50	293	6.1	7.6	2.17	12.90	6.88	999	1041.7	5.6	6.4	999.0	99.0	99.00	
2017	01	01	06	50	314	4.2	5.4	2.24	12.12	7.09	999	1042.1	5.5	6.4	999.0	99.0	99.00	
2017	01	01	07	50	297	4.4	5.2	1.97	13.79	6.64	999	1041.9	5.5	6.4	999.0	99.0	99.00	
2017	01	01	08	50	287	4.6	5.6	2.06	12.90	7.26	999	1041.8	5.5	6.4	999.0	99.0	99.00	
2017	01	01	09	50	304	3.2	4.1	1.82	12.12	6.65	999	1041.5	5.6	6.4	999.0	99.0	99.00	
2017	01	01	10	50	274	2.9	3.9	1.92	12.90	6.88	999	1041.4	5.5	6.4	999.0	99.0	99.00	
2017	01	01	11	50	287	1.3	1.7	1.60	12.90	6.66	999	1041.1	5.4	6.4	999.0	99.0	99.00	
2017	01	01	12	50	260	2.0	3.0	1.67	12.90	6.95	999	1040.7	5.4	6.4	999.0	99.0	99.00	
.
.
2017	12	31	11	50	999	14.5	17.7	5.70	9.09	7.67	134	999.3	5.4	4.9	999.0	99.0	99.00	
2017	12	31	12	50	999	11.8	14.7	6.21	10.00	8.52	125	1001.0	5.5	4.9	999.0	99.0	99.00	
2017	12	31	13	50	999	10.1	12.6	5.52	10.81	8.49	165	1002.0	5.3	4.9	999.0	99.0	99.00	
2017	12	31	14	50	999	8.6	10.5	5.47	10.00	8.61	150	1002.5	5.3	4.9	999.0	99.0	99.00	
2017	12	31	15	50	999	7.1	9.1	5.83	12.90	9.27	216	1002.7	5.2	4.9	999.0	99.0	99.00	
2017	12	31	16	50	999	6.4	8.4	5.82	11.43	9.35	177	1002.7	5.1	4.9	999.0	99.0	99.00	
2017	12	31	17	50	999	5.3	7.4	5.52	11.43	9.04	182	1002.4	5.0	4.9	999.0	99.0	99.00	
2017	12	31	18	50	999	2.8	4.6	5.00	11.43	8.82	179	1002.0	4.8	4.9	999.0	99.0	99.00	
2017	12	31	19	50	999	3.4	5.0	4.87	12.12	8.63	200	1001.6	5.0	4.9	999.0	99.0	99.00	
2017	12	31	20	50	999	3.0	4.4	4.98	12.90	8.57	201	1000.6	4.8	4.9	999.0	99.0	99.00	
2017	12	31	21	50	999	3.8	6.3	4.64	10.00	8.55	150	1000.1	4.8	4.9	999.0	99.0	99.00	
2017	12	31	22	50	999	3.4	5.2	4.40	12.90	8.40	200	998.9	4.7	4.9	999.0	99.0	99.00	

Date	Pres	TIDE
2017 01 01 00 00	106400	103300
2017 01 01 06 00	106400	103280
2017 01 01 12 00	106380	103220
2017 01 01 18 00	106270	103030
2017 01 02 00 00	106150	103030
2017 01 02 06 00	106030	102920
2017 01 02 12 00	105960	102820
2017 01 02 18 00	105980	102850
2017 01 03 00 00	106010	102840
2017 01 03 06 00	106080	102880
2017 01 03 12 00	106100	102960
2017 01 03 18 00	106090	102920
2017 01 04 00 00	106070	102920
2017 01 04 06 00	106040	102810
.	.	.
.	.	.
.	.	.
2017 12 31 12 00	102520	99040
2017 12 31 18 00	102660	99400

Figure 7: An example of matching the data from NDBC (left) and NNRP (right).

- *Classification*: The final datasets will be ready to use as input for ML classifiers, requiring a nominal output attribute, whose specific preparation is detailed in Section 3.5.
- *Regression*: The final datasets will be ready to use as input for regression methods, requiring a real output attribute, whose preparation is also explained in Section 3.5.
- *Direct matching*: In this case the inputs attributes have a direct correspondence with the output attribute, and it is not necessary to perform any additional preparation. Both input and target attributes are synchronised in time, in such a way that the final dataset is not intended for prediction purposes. For example, the final datasets may be used in lost data recovering tasks, in correlation studies, in descriptive analyses, etc.

The following parameters can be specified for the matching process:

- *Flux of energy* [29]: When the F_e is selected, it will be used as output. This attribute is not collected by the buoys, but it can be calculated from two wave parameters: H_s and T_e , which are collected as WVHT and APD attributes, respectively, and were described in Table 1. In this way,

SPAMDA obtains the F_e of each instance using the following equation:

$$F_e = 0.49 \cdot H_s^2 \cdot T_e, \quad (1)$$

where F_e is measured in kilowatts per meter, H_s is measured in meters and T_e is measured in seconds. Note also that F_e is defined in Eq. 1 as an average energy flux (H_s is a kind of average wave height), though for simplicity it will be referred just as flux of energy.

- *Attribute to predict*: Instead of using F_e , researchers can select any of the attributes collected by the buoys as output (e.g. significant wave height, WVHT, wind direction, WDIR, sea level pressure, PRES, etc.). Therefore, they can conduct different studies by selecting one attribute or other.

- *Reanalysis data files*: In order to have a possible better description of the problem under study, more than one reanalysis variable can be considered as input. Remember that these files have to be previously downloaded from the NNRP website [37], which should set the range of dates (temporal properties) and the desired sub-grid (spatial properties, see Fig. 2) for each variable of reanalysis.

In that sense, the reanalysis data files must have the same spatial and temporal properties but related to different variables. SPAMDA simplifies this task by showing the reanalysis data files that are compatibles each other, and checking that the selection made by the researches meets that condition.

- *Buoys attributes*: In addition to the reanalysis variables, the final datasets will also include the selected attributes as inputs (of the intermediate or pre-processed dataset used), providing a possible better characterisation of the problem under study, although it will depend on how correlated the attributes are.
- *Include missing dates*: As above-mentioned, the information collected by a buoy may be incomplete due to measurements not recorded by it. As a

390 consequence, the matching of instances between both sources of information may not be possible (missing dates). In that situation, the researchers can consider two options: 1) discard the instances affected or 2) include them. In the latter case, the final datasets will contain the affected instances, but the measurements of the buoy will be stored as missing values
 395 in WEKA format, denoted as «?».

- *Nearest reanalysis nodes to consider:* As already shown in Fig. 2 (which represents six reanalysis nodes), the reanalysis data files may contain information of several reanalysis nodes. In this way, the researcher can:

- Consider all the reanalysis nodes contained in each file: in this case,
 400 the information provided by each reanalysis node contained in each selected reanalysis data file will be used.
- Consider only some of the reanalysis nodes contained in each file: in this case, only the information of the N closest reanalysis nodes to the buoy will be used (N given by the user). To do that, SPAMDA
 405 uses the *Haversine* equation [47] to calculate the distance from each reanalysis node to the location of the buoy and obtain the closest ones. Haversine equation is also known as the great circle distance and performs calculation from main point to destination point with a trigonometric function using latitude and longitude:

$$\begin{aligned}
 d(p_0, p_j) = & \arccos(\sin(lat_0) \cdot \sin(lat_j) \\
 & \cdot \cos(lon_0 - lon_j) + \cos(lat_0) \\
 & \cdot \cos(lat_j)),
 \end{aligned} \tag{2}$$

410 where p_0 is the buoy geographical location, p_j stands for the location of each reanalysis node, and lat and lon are the latitude and longitude of the points, respectively.

- *Number of final datasets:* Depending on the number of nearest reanalysis nodes to consider, the number of final datasets to create and the content

415

of them can be configured according to the following options:

420

- *One (using weighted mean of the N nearest reanalysis nodes)*: Only one final dataset will be created, which will contain the attributes (the selected one as output and the selected ones as inputs) of the intermediate or pre-processed dataset used, along with a weighted mean of each variable of the reanalysis data used (one per selected reanalysis data file). This weighted mean is obtained by SPAMDA and uses Eq. 2 to obtain the distance from each reanalysis node to the location of the buoy. Once the distances have been calculated they are inverted and normalised as follows:

425

$$w_i = \frac{\sum_{j=1}^N d(p_0, p_j)}{d(p_0, p_i)}, \quad i = 1, \dots, N. \quad (3)$$

430

With these weights, a weighted mean of each variable of reanalysis is obtained for each of the N nodes. Therefore, the closest reanalysis nodes to the localisation of the buoy will provide more information. Considering as example the two nearest reanalysis nodes represented in Fig. 2 and the reanalysis variables air temperature and pressure, the weighted mean of each reanalysis variable will be calculated using the reanalysis nodes $57.5 \text{ N} \times 147.5 \text{ W}$ and $55.0 \text{ N} \times 147.5 \text{ W}$.

435

- *'N' (one per each reanalysis node)*: As many final datasets as the number of nearest N reanalysis nodes configured by the researcher will be created. Therefore, each final dataset will contain the value of each reanalysis variable used of the nearest corresponding reanalysis node, along with the selected attributes of the intermediate or pre-processed dataset used.

440

In this case, and considering as example the four closest reanalysis nodes (see Fig. 2) and the reanalysis variables air temperature and pressure, four final datasets will be created, containing each one the information of both reanalysis variables of the corresponding reanalysis node: $57.5 \text{ N} \times 147.5 \text{ W}$, $55.0 \text{ N} \times 147.5 \text{ W}$, $57.5 \text{ N} \times 150.0$

W and 55.0 N \times 150.0 W, along with the selected attributes of the intermediate or pre-processed dataset used.

Once the matching parameters have been described, for a better understanding of them, Fig. 8 presents an example of the matched information considering the data shown in Fig. 7 and using the following matching configuration¹:

- Variable WVHT as attribute to predict.
- Variable Pres as reanalysis input attribute.
- Variable WSDP as buoy input attribute.
- Not including missing dates.
- Considering the closest reanalysis node.
- Task to be used: *Direct matching*.

Date	Pres	WSPD	WVHT
2017 01 01 00 00	106400	6.4	2.41
2017 01 01 06 00	106400	6.1	2.17
2017 01 01 12 00	106380	1.3	1.60
2017 01 01 18 00	106270	2.4	1.33
2017 01 02 00 00	106150	1.2	0.94
2017 01 02 06 00	106030	3.0	1.42
2017 01 02 12 00	105960	3.1	1.99
2017 01 02 18 00	105980	2.5	2.52
2017 01 03 00 00	106010	3.5	2.03
2017 01 03 06 00	106080	5.1	1.84
2017 01 03 12 00	106100	5.4	1.81
2017 01 03 18 00	106090	5.6	1.60
2017 01 04 00 00	106070	7.0	1.54
2017 01 04 06 00	106040	7.4	1.73
.	.	.	.
.	.	.	.
.	.	.	.
2017 12 31 12 00	102520	14.5	5.70
2017 12 31 18 00	102660	5.3	5.52

Figure 8: An example of the matched information for *Direct matching*.

¹Note that the date is shown just for better understanding, but it will not be included in the final dataset.

455 3.5. Final datasets

Once the matching process has been performed with the desired configuration, it is necessary to prepare the matched information for the desired prediction task (*Regression* or *Classification*), obtaining as a result the final datasets. Remember that *Direct matching*, as it was described in Section 3.4, performs
460 a direct correspondence between the attributes used as inputs and the output one, and it is not necessary to carry out any preparation.

SPAMDA allows the researchers to make such preparation by means of the following options:

- *Prediction horizon* (Classification and Regression): This option indicates
465 the time gap for moving backward the attribute to predict (output attribute). In this way, the input attributes (variables of the buoy and reanalysis data) will be used to predict the output attribute in a specific future time (e.g. +6h, +12h, +18h, +1 day, etc.).

The minimum interval for increasing and decreasing the prediction horizon
470 is 6h (due to reanalysis data temporal resolution) [2], the same interval used when the matching process is carried out. Therefore, for each increment of the prediction horizon, an instance of the dataset is lost (as this future information is not available). As the minimum prediction horizon is 6h, at least one instance will be lost. The relation between the
475 inputs and the attribute to predict will be defined as follows:

$$o_{t+\Delta t} = \phi(\mathbf{b}_t, \mathbf{r}_t), \quad (4)$$

where t represents the time instant to study and Δt the prediction horizon; o is the attribute to be predicted, \mathbf{b}_t is the vector containing the selected NDBC variables and \mathbf{r}_t is the vector containing the selected reanalysis
480 variables. In this way and considering the matched information shown in Fig. 8, WVHT is o , the vector \mathbf{b} contains the variable WSPD and the vector \mathbf{r} contains Pres.

Optionally, the reanalysis variables can be synchronised with the attribute to predict. Given that these variables are estimated by a mathematical model, we can obtain very good future estimations, which can improve the performance of the results. In this case, the relation between the inputs and the attribute to predict would be:

$$o_{t+\Delta t} = \phi(\mathbf{b}_t, \mathbf{r}_{t+\Delta t}). \quad (5)$$

Note that the selected NDBC variables as input cannot be synchronised with the attribute to predict.

For better understanding, considering the matched information shown in Fig. 8, an example of the creation of one *Regression* dataset is shown in Fig. 9. As mentioned earlier, this prediction task requires a real output variable (in this case, WVHT, the last one). The options considered for the preparation of each final dataset are the following:

- Do not synchronise the reanalysis data (see Eq. 4 for the relation between the inputs and the output).
- A prediction horizon of 6h.

Date	Pres	WSPD	WVHT
2017 01 01 00 00	106400	6.4	2.17
2017 01 01 06 00	106400	6.1	1.60
2017 01 01 12 00	106380	1.3	1.33
2017 01 01 18 00	106270	2.4	0.94
2017 01 02 00 00	106150	1.2	1.42
2017 01 02 06 00	106030	3.0	1.99
2017 01 02 12 00	105960	3.1	2.52
2017 01 02 18 00	105980	2.5	2.03
2017 01 03 00 00	106010	3.5	1.84
2017 01 03 06 00	106080	5.1	1.81
2017 01 03 12 00	106100	5.4	1.60
2017 01 03 18 00	106090	5.6	1.54
2017 01 04 00 00	106070	7.0	1.73
2017 01 04 06 00	106040	7.4	1.64
.	.	.	.
.	.	.	.
.	.	.	.
2017 12 31 12 00	102520	14.5	5.52

Figure 9: An example of the creation of a *Regression* dataset, with a prediction horizon of 6h and without synchronisation.

Note that, due to prediction horizon is 6h, the values of WVHT attribute are moved backward one instance (up). As a consequence, the last instance

(2017/12/31 18:00) is lost and is not included in the final dataset. Besides, and because the reanalysis data has not been synchronised, the values of the Pres and WSPD variables are at the same time instant (t in Eq. 4).

Moreover, considering again the matched information shown in Fig. 8, an example of the creation of the same dataset but applying synchronisation (see Eq. 5) is shown in Fig. 10.

Date	Pres	WSPD	WVHT
2017 01 01 00 00	106400	6.4	2.17
2017 01 01 06 00	106380	6.1	1.60
2017 01 01 12 00	106270	1.3	1.33
2017 01 01 18 00	106150	2.4	0.94
2017 01 02 00 00	106030	1.2	1.42
2017 01 02 06 00	105960	3.0	1.99
2017 01 02 12 00	105980	3.1	2.52
2017 01 02 18 00	106010	2.5	2.03
2017 01 03 00 00	106080	3.5	1.84
2017 01 03 06 00	106100	5.1	1.81
2017 01 03 12 00	106090	5.4	1.60
2017 01 03 18 00	106070	5.6	1.54
2017 01 04 00 00	106040	7.0	1.73
2017 01 04 06 00	105950	7.4	1.64
.	.	.	.
.	.	.	.
.	.	.	.
2017 12 31 12 00	102660	14.5	5.52

Figure 10: An example of the creation of a *Regression* dataset, with a prediction horizon of 6h and with synchronisation.

Again, and due to the prediction horizon selected (6h), the values of the WVHT attribute are moved backward one instance (up) and the last instance (2017/12/31 18:00) is not included in the final dataset. But now, the values of the Pres variable are also moved backward one instance (due to the synchronisation). Therefore, in this case, Pres is at the same time instant as the attribute to predict ($t + \Delta t$ in Eq. 5).

- *Thresholds of the output attribute* (Classification): Since the values of the variables collected by the buoys are real numbers, it is necessary to discretise them (convert them from real to nominal values) for the attribute selected as output (attribute to be predicted). SPAMDA allows the researchers to perform this process by defining the necessary classes with their thresholds, which will be used to carry out such discretisation.

Considering again the matched information shown in Fig. 8, an example

of the creation of a *Classification* dataset is shown in Fig. 11. The options considered for the preparation of the final dataset are the following:

- Do not synchronise the reanalysis data.
- A prediction horizon of 6h.
- The thresholds shown in Table 2.

Table 2: Thresholds for the classification example represented in Fig. 11

Class	Description	Inferior [Superior)
Low	Low wave height	0.36	1.5
Average	Average wave height	1.5	2.5
Big	Big wave height	2.5	4.0
Huge	Huge wave height	4.0	9.9

Date	Pres	WSPD	Class_WVHT
2017 01 01 00 00	106400	6.4	Average
2017 01 01 06 00	106400	6.1	Average
2017 01 01 12 00	106380	1.3	Low
2017 01 01 18 00	106270	2.4	Low
2017 01 02 00 00	106150	1.2	Low
2017 01 02 06 00	106030	3.0	Average
2017 01 02 12 00	105960	3.1	Big
2017 01 02 18 00	105980	2.5	Average
2017 01 03 00 00	106010	3.5	Average
2017 01 03 06 00	106080	5.1	Average
2017 01 03 12 00	106100	5.4	Average
2017 01 03 18 00	106090	5.6	Average
2017 01 04 00 00	106070	7.0	Average
2017 01 04 06 00	106040	7.4	Average
.	.	.	.
.	.	.	.
.	.	.	.
2017 12 31 12 00	102520	14.5	Huge

Figure 11: An example of the creation a *Classification* dataset, with a prediction horizon of 6h and without synchronisation.

Note that the attribute to be predicted has been renamed to *Class_WVHT* to show that it is now a nominal variable, because its values have been discretised according to the thresholds. Besides, and due to the 6h prediction horizon, the last instance is lost (2017/12/31 18:00) and the values of the attribute *Class_WVHT* are moved backward one instance (up). As

530 the reanalysis data have not been synchronised, the values of the Pres and
WSPD variables are at the same time instant (t in Eq. 4).

The content of the final datasets, obtained as the result of the preparation
of the matched data, can be visualised to check everything before saving them
on disk. Such preparation can be performed as many times as required and
535 considering the different options in each moment. Although the date will not be
included in the final datasets, it can be shown to properly check the matching.

Finally, it is necessary to define the output configuration to create the final
datasets:

- *Output path file*: Name of the final datasets and folder to save them on
540 disk.
- *Final datasets format*:
 - *ARFF: Attribute-Relation File Format* [34], which is used by WEKA.
SPAMDA allows the researchers to directly open the final datasets in
the Explorer environment of WEKA (in the same context of work),
545 enabling them to choose the most appropriate ML method to tackle
the problem under study.
 - *CSV: Comma-Separated Values*. This format is included in order to
consider other different tasks of software tools.

A text file that summarises the configuration used in matching process and
550 in the preparation of the matched data is also generated. It can be saved and
loaded, enabling the researchers to resume their studies at any other time.

3.6. Manage reanalysis data

As mentioned in Section 2, the reanalysis data files provided by NNRP
contain the estimated values by a mathematical model of one meteorological
555 variable.

In this module (see Fig. 3), SPAMDA includes features for entering new
files and deleting the unnecessary ones. Besides, useful information about the

content of each reanalysis file can be consulted such as name of the file and the reanalysis variable, number of instances and reanalysis nodes, initial and final
560 time, latitude and longitude. All these fields summarise the temporal and spatial properties of the data. Thus, the researcher can quickly and easily identify each reanalysis file entered in SPAMDA.

An example where two reanalysis data files have been entered in SPAMDA is shown in Fig. 12.

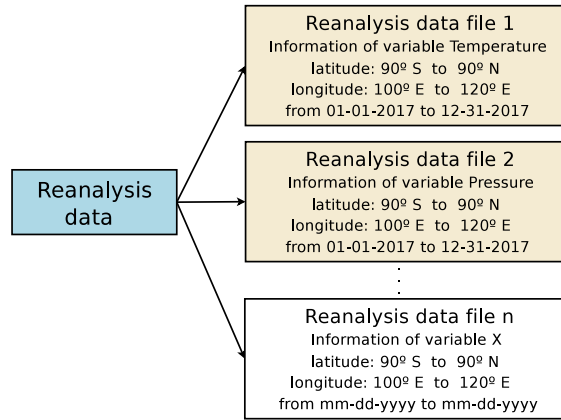


Figure 12: Example of entering two reanalysis data files.

565 3.7. Tools

SPAMDA also contains another module that provides two utilities: one of them is *Dataset converter* used for converting the desired intermediate or pre-processed datasets to ARFF or CSV formats; the other utility can be used for opening ARFF files with WEKA Explorer environment, which is useful for
570 easily checking the results of different configurations of the pre-processing.

4. Case study

Although the software includes a user manual, this section will describe how the application works in a practical approach. To do so, an example showing how to create a fully processed dataset (final dataset), starting from the raw

575 data, will be described. The objective of this final dataset is to be used with ML
algorithms to classify waves in the Gulf of Alaska depending on their height.
The data collected to perform this example is:

1. The measurements obtained from 2013 to 2017 by the buoy with ID 46001,
placed in the Gulf of Alaska, which are provided by NDBC as annual text
580 files. This data is publicly available at the NDBC website.
2. Complementary information collected from reanalysis data containing air
temperature, pressure and two components of wind speed measurements
(South-North and West-East). This information will be collected from the
four closest reanalysis nodes surrounding the geographical location of the
585 buoy. This data is publicly available at the NNRP website and can be
downloaded in NetCDF format.

After gathering the information described above², the researcher can open
SPAMDA.

In Figure 13, the main window is shown. In order to input the reanalysis data
590 which will be used in further steps for creating the final dataset, the researcher
has to select the option *Manage reanalysis data*.

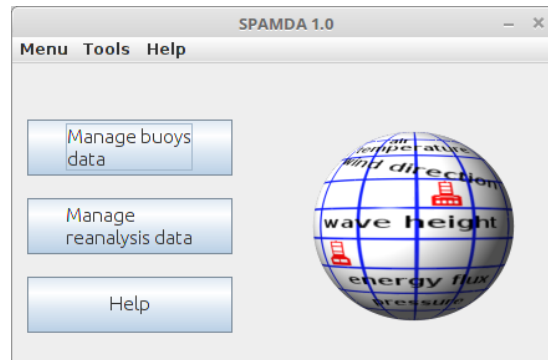


Figure 13: SPAMDA main window.

Then, the window of Figure 14 is shown. Here, using the buttons located at

²Further instructions for downloading this data can be found in the user manual of the application.

the bottom, it is possible to add, delete or consult any data from the different reanalysis files. After the information has been introduced in the application, this window can be closed and the user can go back to the main window to continue entering the information related to the buoy under study.

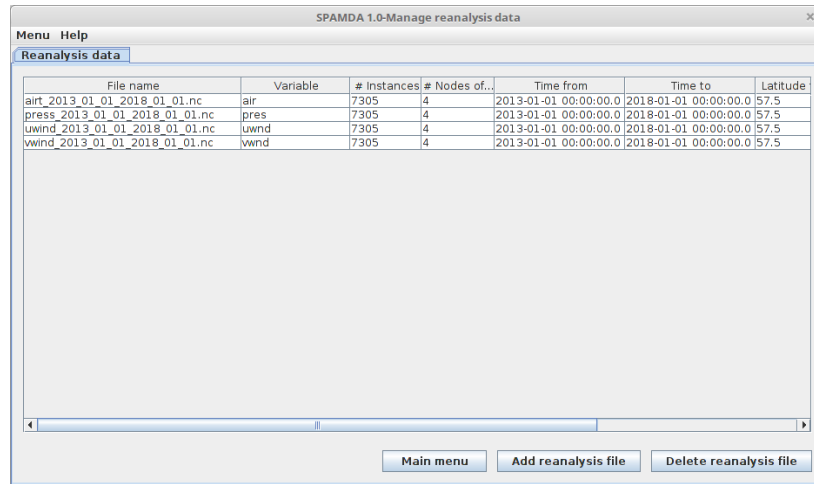


Figure 14: *Manage reanalysis data* window.

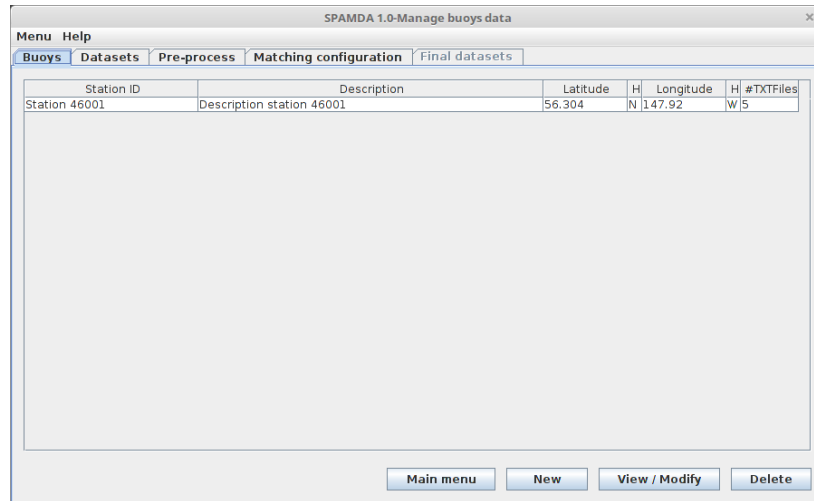


Figure 15: *Buoys* tab.

The researcher has to select *Manage buoys data* to open the window shown

in Figure 15. In this tab, the researcher can consult, modify, add or delete different data related to the buoy. In order to enter such data, click on the *New* button, and then the window shown in Figure 16 pops up.

Figure 16: *Entering a new buoy window.*

Here the information about the buoy has to be included: the *Station ID*, its description, geographical localisation and the corresponding annual text files. In this case, the files containing the data from year 2013 to 2017 are inserted by clicking on the *Add file* button. Once the data has been introduced, it is necessary to click on the *Save* button to insert the buoy in SPAMDA database. After that, the window can be closed.

To create the intermediate dataset, the researcher has to double-click on the buoy under study or click on the *Datasets* tab to switch to the corresponding view (see Figure 17). In this view, the researcher can delete or consult a summary of each intermediate or pre-processed dataset by selecting it from the corresponding list. It can also create new ones. To proceed with the creation of the intermediate dataset, the user clicks on the *New* button, and the view shown in Figure 18 appears.

Here the researcher can select the annual text files to be included in the

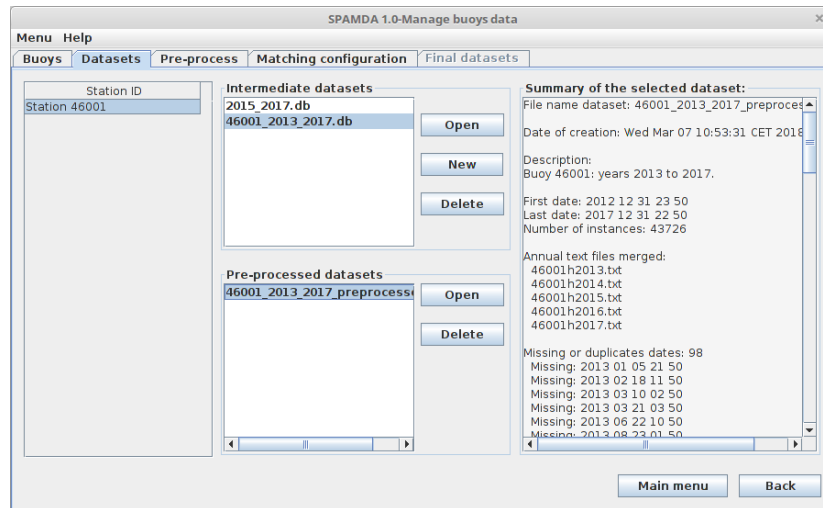


Figure 17: *Datasets* tab.

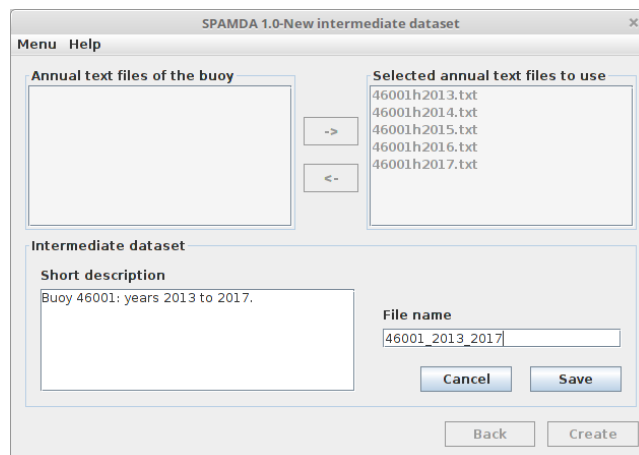


Figure 18: *New intermediate dataset* view.

intermediate dataset. In this case, all the files introduced before, which correspond to the buoy under study, are selected. When the file selection is finished, *Create* button has to be clicked in order to introduce the description and the file name of the current intermediate dataset, and then, with the *Save* button, the creation process starts, showing the status of the process during it.

After that, in order to prepare the intermediate dataset, the dataset is se-

lected, and then the button *Open* is clicked to jump to the tab *Pre-process* (shown in Figure 19). In this tab, relevant statistical information about the selected dataset is shown, and also the content of the dataset can be consulted, providing the researcher the capacity to evaluate the pre-processing being performed.

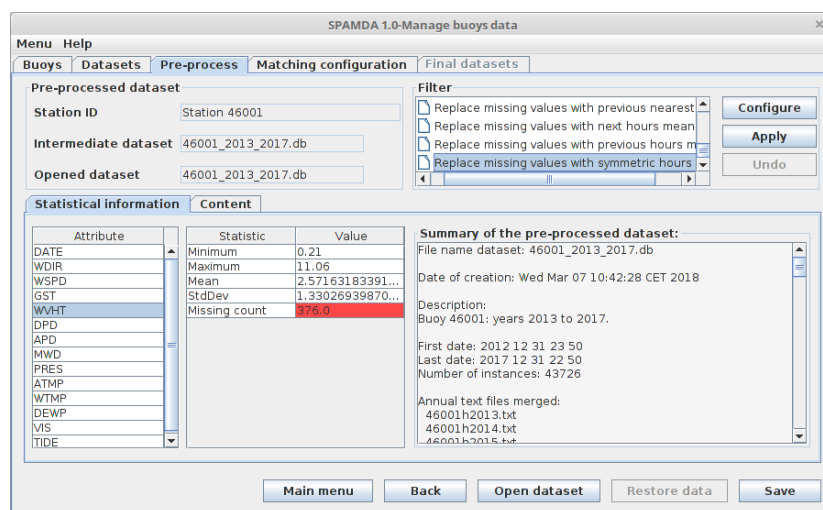


Figure 19: *Pre-process* tab.

Here the researcher can apply (and configure) the necessary filters (explained in Section 3.3) to the selected dataset, and, in the bottom part, the main statistics of the dataset are displayed, which can be used to observe the changes produced when applying a filter. As mentioned earlier, this case study is focused on classifying waves considering their height, so any missing data from wave height (376 values) and the remaining attributes are recovered, using the filter *Replace missing values with symmetric 3 hours mean*. Furthermore, the attributes MWD, DEWP, VIS and TIDE are removed from the dataset by applying the filter *RemoveByName*, since the first two had more than 92% of missing data and the last two 100%. After finishing the pre-processing of the dataset, the researcher can click on the *Save* button, to introduce the description and file name for the current pre-processed dataset.

At this point, the researcher has registered the buoy in SPAMDA, then entered its raw data and selected the required data for the problem (intermediate dataset). Finally, the data has been pre-processed in order to be ready for its future use in ML algorithms. In order to achieve a more accurate description of the problem under study, a matching process can be carried out to merge the processed data from NDBC with the reanalysis data (also entered previously) from NNRP. The next step is to click on the *Matching configuration* tab, to open the view shown in Figure 20.

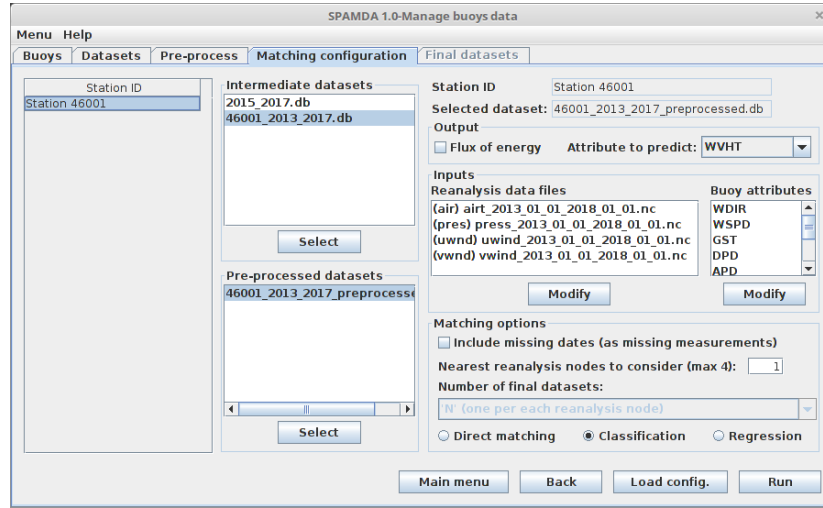


Figure 20: *Matching configuration* tab.

In this view, the researcher can customise (or load) the parameters of the matching process according to their needs and select the prediction task (described in Section 3.4) that the final dataset will be used for. In this example, the following parameters were selected:

- Attribute to predict: WVHT.
- Reanalysis data: Air, pressure, u-wind and v-wind.
- Buoy attributes to be used as inputs: WDIR, WSPD, GST, DPD, APD, PRES, ATMP and WTMP.

- Reanalysis nodes to consider: 1.
- Number of final datasets: In this example that option is disabled, because only one reanalysis node is considered.
- Prediction task: Classification.

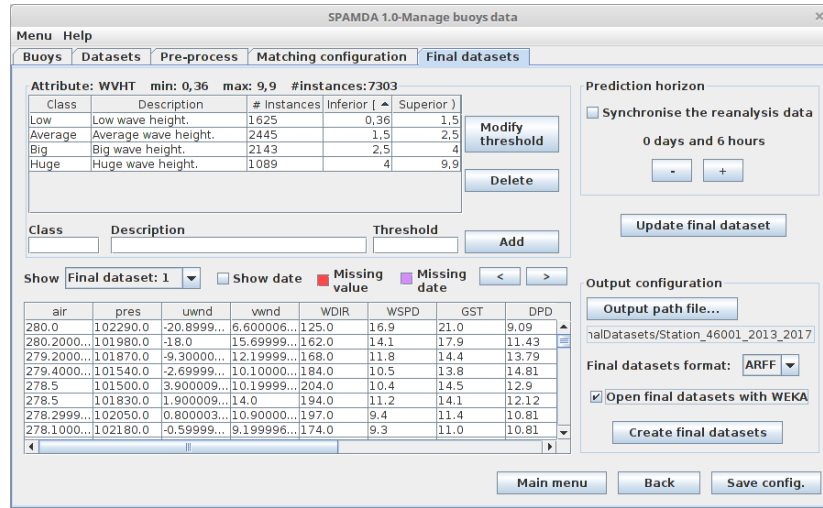


Figure 21: *Final datasets* tab.

After configuring the matching process, the researcher can click on the [Run](#) button to jump to the view shown in Figure 21 and proceed to define the final dataset structure according to the selected prediction task. Given that, in the previous window, *Classification* was selected, the researcher can now add, modify or delete the thresholds (usually defined by an expert) for discretising the output variable. After this, the next step is to set the time horizon desired and also to activate (if desired) the synchronisation (in time) of reanalysis variables with the output, as explained in Section 3.5. Then the researcher can click on the [Update final dataset](#) button to see the content shown in the bottom left corner. Finally, after checking that everything is correct, the last step would be to select the name (and path) of the dataset file and its output format and click on the [Create final datasets](#) button. For this case study, the following configuration was applied:

- Thresholds: see Table 2.
- Prediction horizon: 6 hours
- Synchronisation: Disabled

At this point the final dataset would be created and stored in the computer
 675 of the researcher. Also, there is an option to open the dataset with WEKA (after
 creating it) in order to perform a first classification approach or a preliminary
 study of the data structure, as shown in Fig. 22.

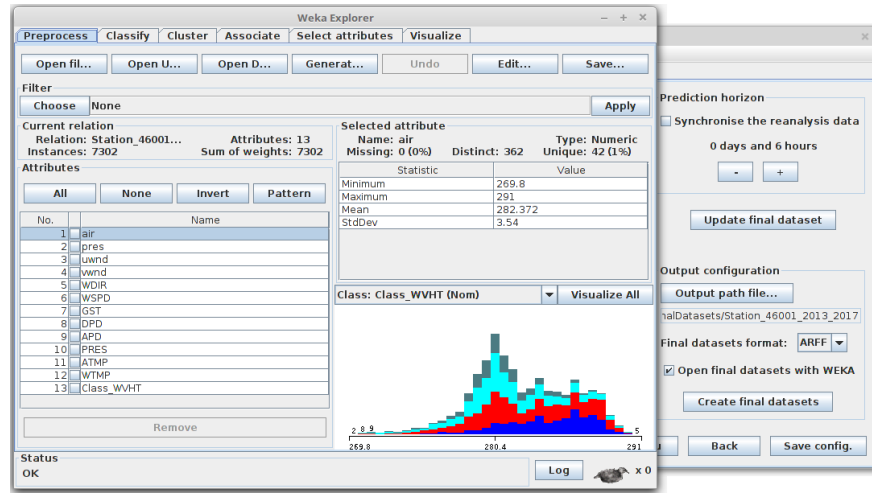


Figure 22: Final dataset opened with the environment Explorer of WEKA.

5. Conclusions

A software tool for creating datasets using meteorological data from NDBC
 680 and NNRP has been presented. These datasets will be ready to use as input for
 ML techniques in prediction tasks (classification or regression). As a result, the
 researchers will benefit from a great support when carrying out their oceanic and
 atmospheric studies, related to energy and environmental modelling. Moreover,
 given that SPAMDA simplifies all the intermediate steps involved in the creation
 685 of datasets (such as entering the meteorological information, managing with

the incomplete data, pre-processing tasks, the customisable matching process to merge the data and the preparation of the datasets according to the ML technique to use), it avoids errors and reduces the time needed. In this way, the researchers will be able to have more in-depth analysis, which could result
690 in more complete conclusions about the issue under study.

In order to improve SPAMDA, some future work could be focused on new functional modules for managing meteorological data from others sources, so that the developed tool can be extended to any other research, new pre-processing functionalities such as filters to analyse the correlation between attributes or new
695 functional modules for recovering missing values using nearby buoys data [48]. Furthermore, the developed software could manage other sources of reanalysis data, and new output formats for the datasets which could be used as input by other tools for ML such as KEEL (*Knowledge Extraction based on Evolutionary Learning*) [49].

700 Acknowledgments

This work has been partially subsidised by the projects TIN2017-85887-C2-1-P and TIN2017-90567-REDT of the Spanish Ministry of Economy and Competitiveness (MINECO), and FEDER funds of the European Union. We also thank to NVIDIA Corporation for the transfer of computational resources
705 for research works.

The authors also thank to NOAA/OAR/ESRL PSD, Boulder, Colorado, USA for the NCEP Reanalysis data provided from their Web site at <https://www.esrl.noaa.gov/psd/>, to NOAA/NDBC by its data that were collected and made freely available, to University of Waikato for the Weka (Waikato
710 Environment for Knowledge Analysis) software tool, to University Corporation for Atmospheric Research/Unidata for the NetCDF (network Common Data Form) Java library and to QOS.ch for the SLF4J (Simple Logging Facade for Java) library.

References

- 715 [1] V. Laface, F. Arena, C. G. Soares, Directional analysis of sea storms, *Ocean Engineering* 107 (Supplement C) (2015) 45–53. doi:10.1016/j.oceaneng.2015.07.027.
- [2] M. Dorado-Moreno, L. Cornejo-Bueno, P. Gutiérrez, L. Prieto, C. Hervás-Martínez, S. Salcedo-Sanz, Robust estimation of wind power ramp events
720 with reservoir computing, *Renewable Energy* 111 (2017) 428–437. doi:10.1016/j.renene.2017.04.016.
- [3] M. J. Dvorak, C. L. Archer, M. Z. Jacobson, California offshore wind energy potential, *Renewable Energy* 35 (6) (2010) 1244–1254. doi:10.1016/j.renene.2009.11.022.
- 725 [4] H. L. Fuchs, G. P. Gerbi, Seascape-level variation in turbulence- and wave-generated hydrodynamic signals experienced by plankton, *Progress in Oceanography* 141 (Supplement C) (2016) 109–129. doi:10.1016/j.pocean.2015.12.010.
- [5] V. d. P. R. da Silva, R. A. e Silva, E. P. Cavalcanti, C. C. Braga, P. V. de Azevedo, V. P. Singh, E. R. R. Pereira, Trends in solar radiation in
730 NCEP/NCAR database and measurements in northeastern Brazil, *Solar Energy* 84 (10) (2010) 1852–1862. doi:10.1016/j.solener.2010.07.011.
- [6] B. Gouldby, F. J. Méndez, Y. Guanche, A. Rueda, R. Mínguez, A methodology for deriving extreme nearshore sea conditions for structural design
735 and flood risk analysis, *Coastal Engineering* 88 (Supplement C) (2014) 15–26. doi:10.1016/j.coastaleng.2014.01.012.
- [7] C. E. C. Nogueira, M. L. Vidotto, F. Toniazzi, G. Debastiani, Software for designing solar water heating systems, *Renewable and Sustainable Energy Reviews* 58 (Supplement C) (2016) 361–375. doi:10.1016/j.rser.2015.
740 12.346.

- [8] C. K. Lo, Y. S. Lim, F. A. Rahman, New integrated simulation tool for the optimum design of bifacial solar panel with reflectors on a specific site, *Renewable Energy* 81 (Supplement C) (2015) 293–307. doi:10.1016/j.renene.2015.03.047.
- 745 [9] T. H. Nguyen, A. Prinz, T. Friisø, R. Nossun, I. Tyapin, A framework for data integration of offshore wind farms, *Renewable Energy* 60 (Supplement C) (2013) 150–161. doi:10.1016/j.renene.2013.05.002.
- [10] J. E. Gutierrez, B. Zamora, J. García, M. R. Peyrau, Tool development based on FAST for performing design optimization of offshore wind turbines: FASTLognoter, *Renewable Energy* 55 (Supplement C) (2013) 69–78.
750 doi:10.1016/j.renene.2012.12.026.
- [11] N. Baker, M. C. Guedes, N. Shaikh, L. Calixto, R. Aguiar, The LT-Portugal software: A design tool for architects, *Renewable Energy* 49 (Supplement C) (2013) 156–160, selected papers from World Renewable Energy Congress - XI. doi:10.1016/j.renene.2012.01.041.
755
- [12] A. L. A. Raabe, A. H. d. F. Klein, M. González, R. Medina, MEPBAY and SMC: Software tools to support different operational levels of headland-bay beach in coastal engineering projects, *Coastal Engineering* 57 (2) (2010) 213–226, hydrodynamics and Applications of Headland-Bay Beaches. doi:10.1016/j.coastaleng.2009.10.008.
760
- [13] K. Skeiker, B. A. Ghani, A software tool for the creation of a typical meteorological year, *Renewable Energy* 34 (3) (2009) 544–554. doi:10.1016/j.renene.2008.05.046.
- [14] M. Zeyringer, B. Fais, I. Keppo, J. Price, The potential of marine energy technologies in the UK – Evaluation from a systems perspective, *Renewable Energy* 115 (Supplement C) (2018) 1281–1293. doi:10.1016/j.renene.2017.07.092.
765

- [15] M. Bhattacharya, S. A. Churchill, S. R. Paramati, The dynamic impact of renewable energy and institutions on economic output and CO2 emissions across regions, *Renewable Energy* 111 (Supplement C) (2017) 157–167. doi:10.1016/j.renene.2017.03.102.
- [16] A. F. d. O. Falcão, Wave energy utilization: A review of the technologies, *Renewable and Sustainable Energy Reviews* 14 (3) (2010) 899–918. doi:10.1016/j.rser.2009.11.003.
- [17] M. K. Ochi, *Ocean Waves: The Stochastic Approach*, Cambridge Ocean Technology Series, Cambridge University Press, 1998. doi:10.1017/CB09780511529559.
- [18] S. Crowley, R. Porter, D. J. Taunton, P. A. Wilson, Modelling of the WITT wave energy converter, *Renewable Energy* 115 (Supplement C) (2018) 159–174. doi:10.1016/j.renene.2017.08.004.
- [19] O. Abdelkhalik, R. Robinett, S. Zou, G. Bacelli, R. Coe, D. Bull, D. Wilson, U. Korde, On the control design of wave energy converters with wave prediction, *Journal of Ocean Engineering and Marine Energy* 2 (4) (2016) 473–483. doi:10.1007/s40722-016-0048-4.
- [20] J. V. Ringwood, G. Bacelli, F. Fusco, Energy-Maximizing Control of Wave-Energy Converters: The Development of Control System Technology to Optimize Their Operation, *IEEE Control Systems* 34 (5) (2014) 30–55. doi:10.1109/MCS.2014.2333253.
- [21] K. Chatziioannou, V. Katsardi, A. Koukouselis, E. Mistakidis, The effect of nonlinear wave-structure and soil-structure interactions in the design of an offshore structure, *Marine Structures* 52 (Supplement C) (2017) 126–152. doi:10.1016/j.marstruc.2016.11.003.
- [22] Y. Dalgic, I. Lazakis, I. Dinwoodie, D. McMillan, M. Revie, Advanced logistics planning for offshore wind farm operation and maintenance ac-

- 795 activities, *Ocean Engineering* 101 (Supplement C) (2015) 211–226. doi:
 10.1016/j.oceaneng.2015.04.040.
- [23] E. Alpaydin, *Introduction to Machine Learning (Adaptive Computation and Machine Learning)*, The MIT Press, 2004.
- [24] C. M. Bishop, *Pattern Recognition and Machine Learning (Information Sci-*
 800 *ence and Statistics)*, Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2006.
- [25] *Large-Scale Machine Learning in the Earth Sciences*, Chapman & Hall/CRC Data Mining and Knowledge Discovery Series.
- [26] M. A. H. Eibe Frank, I. H. Witten, *The WEKA Workbench. Online Ap-*
 805 *pendix for Data Mining: Practical Machine Learning Tools and Techniques* (2016).
- [27] A. M. Durán-Rosal, J. C. Fernández, P. A. Gutiérrez, C. Hervás-Martínez, Detection and prediction of segments containing extreme significant wave heights, *Ocean Engineering* 142 (Supplement C) (2017) 268–279. doi:
 810 10.1016/j.oceaneng.2017.07.009.
- [28] N. K. kumar, R. Savitha, A. A. Mamun, Regional ocean wave height prediction using sequential learning neural networks, *Ocean Engineering* 129 (2017) 605–612. doi:10.1016/j.oceaneng.2016.10.033.
- [29] J. C. Fernández, S. Salcedo-Sanz, P. A. Gutiérrez, E. Alexandre, C. Hervás-Martínez, Significant wave height and energy flux range forecast with machine learning classifiers, *Engineering Applications of Artificial Intelligence* 43 (Supplement C) (2015) 44–53. doi:10.1016/j.engappai.2015.03.
 815 012.
- [30] J. Adams, S. Flora, Correlating seabird movements with ocean winds: linking satellite telemetry with ocean scatterometry, *Marine Biology* 157 (4) (2010) 915–929. doi:10.1007/s00227-009-1367-y.
- 820

- [31] National Data Buoy Center, National Oceanic and Atmospheric Administration of the USA (NOAA), <http://www.ndbc.noaa.gov/>, (Accessed 17 November 2017).
- 825 [32] E. Kalnay, M. Kanamitsu, R. Kistler, W. Collins, D. Deaven, L. Gandin, M. Iredell, S. Saha, G. White, J. Woollen, Y. Zhu, A. Leetmaa, R. Reynolds, M. Chelliah, W. Ebisuzaki, W. Higgins, J. Janowiak, K. C. Mo, C. Ropelewski, J. Wang, R. Jenne, D. Joseph, The NCEP/NCAR 40-Year Reanalysis Project, *Bulletin of the American Meteorological Society* 77 (3) (1996) 437–471. doi:10.1175/1520-0477(1996)077<0437:TNYP>2.0.CO;2.
- 830 [33] R. Kistler, W. Collins, S. Saha, G. White, J. Woollen, E. Kalnay, M. Chelliah, W. Ebisuzaki, M. Kanamitsu, V. Kousky, H. van den Dool, R. Jenne, M. Fiorino, The NCEP–NCAR 50–Year Reanalysis: Monthly Means CD–ROM and Documentation, *Bulletin of the American Meteorological Society* 82 (2) (2001) 247–267. doi:10.1175/1520-0477(2001)082<0247:TNYP>2.3.CO;2.
- 835 [34] The WEKA Data Mining Software: Attribute-Relation File Format (ARFF), <http://https://www.cs.waikato.ac.nz/ml/weka/arff.html>, (Accessed 18 December 2017).
- 840 [35] National Data Buoy Center, NDBC - Historical NDBC Data, http://www.ndbc.noaa.gov/historical_data.shtml, (Accessed 15 January 2018).
- [36] National Data Buoy Center, NDBC - Important NDBC Web Site Changes, <http://www.ndbc.noaa.gov/mods.shtml>, (Accessed 15 January 2018).
- 845 [37] NOAA/OAR/ESRL PSD, ESRL : PSD : NCEP/NCAR Reanalysis 1, <https://www.esrl.noaa.gov/psd/data/gridded/data.ncep.reanalysis.html>, (Accessed 15 January 2018).
- [38] Unidata, Network Common Data Form (NetCDF) version 4.6.10 [software].

- Boulder, CO: UCAR/Unidata., <https://doi.org/10.5065/D6H70CW6>
850 (2017).
- [39] The WEKA Data Mining Software: Normalize, <http://weka.sourceforge.net/doc.stable-3-8/weka/filters/unsupervised/attribute/Normalize.html>, (Accessed 04 December 2017).
- [40] The WEKA Data Mining Software: Remove, <http://weka.sourceforge.net/doc.stable-3-8/weka/filters/unsupervised/attribute/Remove.html>, (Accessed 15 December 2017).
855
- [41] The WEKA Data Mining Software: RemoveByName, <http://weka.sourceforge.net/doc.stable-3-8/weka/filters/unsupervised/attribute/RemoveByName.html>, (Accessed 15 December 2017).
- 860 [42] The WEKA Data Mining Software: ReplaceMissingValues, <http://weka.sourceforge.net/doc.stable-3-8/weka/filters/unsupervised/attribute/ReplaceMissingValues.html>, (Accessed 15 December 2017).
- [43] The WEKA Data Mining Software: ReplaceMissingWithUserConstant, <http://weka.sourceforge.net/doc.stable-3-8/weka/filters/unsupervised/attribute/ReplaceMissingWithUserConstant.html>,
865 (Accessed 15 December 2017).
- [44] The WEKA Data Mining Software: RemoveDuplicates, <http://weka.sourceforge.net/doc.stable-3-8/weka/filters/unsupervised/instance/RemoveDuplicates.html>, (Accessed 04 December 2017).
870
- [45] The WEKA Data Mining Software: RemoveWithValues, <http://weka.sourceforge.net/doc.stable-3-8/weka/filters/unsupervised/instance/RemoveWithValues.html>, (Accessed 15 December 2017).
- [46] The WEKA Data Mining Software: SubsetByExpression,
875 <http://weka.sourceforge.net/doc.stable-3-8/weka/filters/>

`unsupervised/instance/SubsetByExpression.html`, (Accessed 15 December 2017).

- [47] M. J. de Smith, M. F. Goodchild, P. A. Longley, Geospatial Analysis: A Comprehensive Guide to Principles, Techniques and Software Tools, 3rd Edition, Matador, 2009.
- [48] A. M. Durán-Rosal, C. Hervás-Martínez, A. J. Tallón-Ballesteros, A. C. Martínez-Estudillo, S. Salcedo-Sanz, Massive missing data reconstruction in ocean buoys with evolutionary product unit neural networks, *Ocean Engineering* 117 (2016) 292–301, jCR(2016): 1.894 Position: 2/14 (Q1) Category: ENGINEERING, MARINE. doi:10.1016/j.oceaneng.2016.03.053.
- [49] J. Alcalá-Fdez, L. Sánchez, S. García, M. J. del Jesús, S. Ventura, J. M. G. i Guiu, J. M. Otero, C. Romero, J. Bacardit, V. M. R. Santos, J. C. Fernández, F. Herrera, KEEL: a software tool to assess evolutionary algorithms for data mining problems, *Soft Comput.* 13 (2009) 307–318. doi:10.1007/s00500-008-0323-y.
- [50] National Data Buoy Center, NDBC - Measurement Descriptions and Units, <https://www.ndbc.noaa.gov/measdes.shtml>, (Accessed 15 January 2018).

Appendix A. Managing incomplete data

In this appendix, we describe how SPAMDA deals with incomplete data when creating intermediate datasets and performing the matching process.

The measurements collected by the buoys may be incomplete or recorded at a different time than the expected one, due to the weather conditions in which the buoys have to operate. To illustrate this casuistry, the following examples are shown in Fig. 23:

- In the instance marked with a), the measurement of 17:50 was collected at 17:45, 5 minutes earlier.
- In the instance marked with b), the measurement of 23:50 was collected at 23:30, 20 minutes earlier.
- In the instance marked with c), the measurement of 05:50 is duplicated.
- In the instance marked with d), the measurement of 11:50 is missing (missing date or instance).
- In the instance marked with e), the measurement of 17:50 and 18:50 are missing (missing dates or instances).
- Missing values highlighted in red colour.

	#YY	MM	DD	hh	mm	WDIR	WSPD	GST	WVHT	DPD	APD	MWD	PRES	ATMP	WTMP	DEWP	VIS	TIDE
	#yr	mo	dy	hr	mn	degT	m/s	m/s	m	sec	sec	degT	hPa	degC	degC	degC	mi	ft
	2016	12	31	23	50	279	6.4	7.3	2.41	12.90	6.50	999	1041.3	5.6	6.4	999.0	99.0	99.00
	2017	01	01	00	50	291	6.3	7.3	2.13	7.14	6.08	999	1041.1	5.5	6.4	999.0	99.0	99.00

a) ->	2017	01	04	16	50	308	5.5	6.7	1.28	10.00	5.96	999	1037.2	3.8	6.4	999.0	99.0	99.00
	2017	01	04	17	45	297	4.5	5.5	1.24	10.00	5.62	999	1037.3	5.3	6.3	999.0	99.0	99.00
	2017	01	04	18	50	318	5.8	7.1	1.37	10.00	6.04	999	1037.7	5.6	6.3	999.0	99.0	99.00

b) ->	2017	01	04	22	50	329	5.3	6.4	1.16	8.33	5.22	999	1037.0	6.0	6.4	999.0	99.0	99.00
	2017	01	04	23	30	327	5.1	6.2	0.85	5.00	4.13	999	1036.8	5.9	6.4	999.0	99.0	99.00
	2017	01	05	00	50	334	5.6	6.6	1.07	10.81	4.73	999	1036.8	6.0	6.4	999.0	99.0	99.00

	2017	01	05	04	50	321	6.1	7.1	1.33	7.14	5.04	999	1036.6	5.8	6.4	999.0	99.0	99.00
	2017	01	05	05	50	308	6.9	8.0	1.66	7.69	5.47	999	1036.5	5.9	6.4	999.0	99.0	99.00
c) ->	2017	01	05	05	50	308	6.9	8.0	1.66	7.69	5.47	999	1036.5	5.9	6.4	999.0	99.0	99.00

	2017	01	05	10	50	329	10.5	13.1	2.55	8.33	5.44	999	1036.0	5.8	6.3	999.0	99.0	99.00
d) ->	2017	01	05	11	50	(missing date)												
	2017	01	05	12	50	337	11.3	14.1	2.77	7.69	5.69	999	1035.8	5.7	6.3	999.0	99.0	99.00

	2017	01	05	15	50	344	12.3	14.3	3.07	9.09	5.99	999	1034.9	5.7	6.3	999.0	99.0	99.00
	2017	01	05	16	50	345	11.4	14.6	3.16	10.00	6.27	999	1034.9	5.7	6.3	999.0	99.0	99.00
e) ->	2017	01	05	17	50	(missing date)												
	2017	01	05	18	50	(missing date)												
	2017	01	05	19	50	331	11.5	14.2	2.98	9.09	5.84	999	1034.9	5.5	6.3	999.0	99.0	99.00

	2017	12	31	17	50	999	5.3	7.4	5.52	11.43	9.04	182	1002.4	5.0	4.9	999.0	99.0	99.00
	2017	12	31	18	50	999	2.8	4.6	5.00	11.43	8.82	179	1002.0	4.8	4.9	999.0	99.0	99.00
	2017	12	31	19	50	999	3.4	5.0	4.87	12.12	8.63	200	1001.6	5.0	4.9	999.0	99.0	99.00
	2017	12	31	20	50	999	3.0	4.4	4.98	12.90	8.57	201	1000.6	4.8	4.9	999.0	99.0	99.00
	2017	12	31	21	50	999	3.8	6.3	4.64	10.00	8.55	150	1000.1	4.8	4.9	999.0	99.0	99.00
	2017	12	31	22	50	999	3.4	5.2	4.40	12.90	8.40	200	998.9	4.7	4.9	999.0	99.0	99.00

Figure 23: A fragment of an annual text file with different missing value examples.

SPAMDA has been designed to tackle these situations, and it informs the researchers of any incidence found while reading the annual text files for creating the intermediate datasets. For the case of measurements that were recorded at
 915 a different time than expected, it has been established a time gap of 6 minutes (10% of an hour). Therefore, if the time difference exceeds such value the date will be considered as an unexpected.

Fig. 24 shows the status of the creation of an intermediate dataset with the information of Fig 23. Note that the instance marked with a) has not been
 920 informed by SPAMDA as an unexpected date because its time difference is less than 6 minutes. Depending on the affected attribute, NDBC uses a specific value [50] to indicate the presence of lost data (e.g. 99 for VIS and TIDE attributes, 999 for DEWP, MWD and WDIR, etc.). SPAMDA interprets these specific values and, after creating the intermediate dataset, the researchers can check if
 925 it contains missing values by visualising its statistical information or content. Remember that SPAMDA provides several filters for recovering missing data, which were described in subsection 3.2.

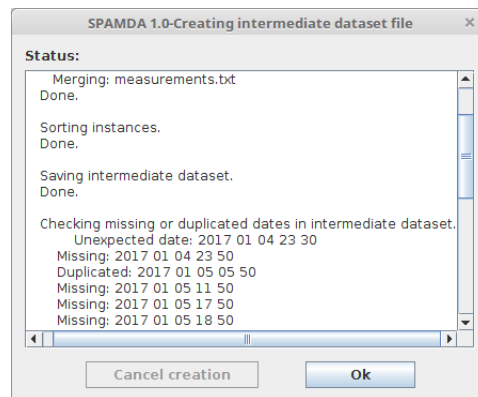


Figure 24: Status of the creation of the intermediate dataset for the example of Fig 23.

SPAMDA takes into account this casuistry when carrying out the matching process. An example is given in Fig. 25. As above-mentioned, the matching
 930 process is performed with the nearest measurement (previous or next) within a maximum of 60 minutes of difference. However, in the instance marked with

