

Surrogate-assisted and filter-based multi-objective evolutionary feature selection for deep learning

Raquel Espinosa*, Fernando Jiménez† and José Palma† *Senior Member, IEEE*

*International Doctorate School. University of Murcia, 30100, Murcia, Spain

†Department of Information Engineering and Communications, University of Murcia, Spain

Abstract—Feature selection for deep learning prediction models is a difficult topic for researchers to tackle. Most of the approaches proposed in the literature consist of embedded methods through the use of hidden layers added to the neural network architecture that modify the weights of the units associated with each input attribute so that the worst attributes have less weight in the learning process. Other approaches used for deep learning are filter methods, which are independent of the learning algorithm, which can limit the precision of the prediction model. Wrapper methods are impractical with deep learning due to their high computational cost. In this paper, we propose new attribute subset evaluation feature selection methods for deep learning of the wrapper, filter and wrapper-filter hybrid types, where multi-objective and many-objective evolutionary algorithms are used as search strategies. A novel surrogate-assisted approach is used to reduce the high computational cost of the wrapper-type objective function, while the filter-type objective functions are based on correlation and an adaptation of the reliefF algorithm. The proposed techniques have been applied in a time series forecasting problem of air quality in the Spanish south-east and an indoor temperature forecasting problem in a domestic house, with promising results compared to other feature selection techniques used in the literature.

Index Terms—Feature selection, deep learning, surrogate-assisted, multi-objective evolutionary algorithms, time series forecasting, air quality, indoor temperature.

I. INTRODUCTION

THE increase in the amount and complexity of available data leads to an increase in the dimensionality of the information. *Feature selection* (FS) [1] is a process that reduces the dimensionality of the input data, as well as the complexity of the models created from the input data. However, very few works focus on FS in deep learning. Most of them consist of embedded methods that integrate FS into the training process of deep learning models. Other FS methods widely used with deep learning are filter-type. They are fast methods, but independent of the learning algorithm, so they are good in general but not so accurate for a given learning algorithm. On the other hand, wrapper methods are very accurate but very expensive, especially when the learning algorithm is slow and dataset contains a large number of input attributes. For this reason, little research has been carried out. There are also wrapper-filter hybrids and ensembles for FS.

Therefore, we propose different multi-objective optimization models for FS both wrapper and filter type, and also hybrid wrapper-filter models. The proposed optimization models contain 2, 3 or 4 objectives, which are solved with *multi-objective evolutionary algorithms* (MOEAs) [2]. For the eval-

uation of the objective functions of the wrapper-type, the *root mean squared error* (RMSE) of a surrogate for a *long short-term memory* (LSTM) [3] artificial recurrent neural network is used. Filter-type objectives are based on correlation metrics and an adaptation of the reliefF method for evaluating subsets of attributes. All the proposed optimization models contain an objective for minimizing the number of attributes.

The proposed methods have been applied to a time series forecasting problem of air quality. In particular, to the prediction of the concentration of nitrogen dioxide (NO_2) in the town of La Aljorra, Region of Murcia, Spain. NO_2 is a noxious gas that is produced in combustion processes such as vehicle engines or some industries. In addition, NO_2 interacting with other chemicals in the air like water or oxygen can cause acid rain. Prolonged exposure to this chemical compound can lead to respiratory [4] and cardiovascular [5] diseases. Additionally, a problem of indoor temperature forecasting in a domestic house in Valencia (Spain) has also been analyzed to strengthen the conclusions.

The main contributions of this work are, in summary, the following:

- We propose new multi-objective FS methods of wrapper, filter and hybrid types based on LSTM recurrent neural networks, correlation, the reliefF algorithm and the minimization of the number of attributes.
- We propose a novel surrogate-assisted MOEA to solve the proposed wrapper and hybrid FS methods.
- We apply the proposed FS methods to the forecast of time series for air quality in the city of La Aljorra (Spain), and in a problem of indoor temperature forecasting in a domestic house in Valencia (Spain), using a future prediction horizon to h -step ahead.
- We propose a decision-making mechanism for choosing the non-dominated solution from the Pareto front, specific for time series forecasting problems.
- We propose a methodology to compare multi-objective optimization models for FS, time series forecasting models, and multi-objective and many-objective evolutionary algorithms, based on non-parametric statistical tests on appropriate performance metrics and multi-criteria aggregation functions at h -step ahead prediction horizons. Our models are compared with the models obtained using other FS techniques used in literature, such as linear regression based wrapper FS methods, random forest based wrapper FS methods, correlation based feature ranking methods, reliefF, and embedded FS methods.

- The proposed method reduces the run time of the FS process, allows dealing with high dimensional problems and reduces the complexity of the prediction models. This reduction in computing time implies a reduction in the carbon footprint, helping to deal with climate change. From a social and economic point of view, our proposal is in line with the objectives of the *Paris Agreement* on climate change (<https://www.un.org/en/climatechange/paris-agreement>), the *Sustainable Development Goals* (<https://sdgs.un.org/>) and the *European Green Deal* (<https://ec.europa.eu/info/strategy/priorities-2019-2024/european-green-deal>). On the other hand, reducing the number of attributes in prediction models contributes to *Explainable Artificial Intelligence* (XIA) [6], which is currently in great social demand and is part of the objectives of, among others, the *White Paper on Artificial Intelligence* of the European Commission (<https://ec.europa.eu/info/publications/white-paper-artificial-intelligence-european-approach-excellence-and-trust>) and the *Executive Order 13960 on Promoting the Use of Trustworthy AI in the Federal Government* of USA (<https://www.federalregister.gov/documents/2020/12/08/2020-27065/promoting-the-use-of-trustworthy-artificial-intelligence-in-the-federal-government>).

The rest of the article is organized as follows. Section 2 presents some background for FS, surrogate-assisted evolutionary computation and air quality; Section 3 describes the materials and methods used in this work; Section 4 presents the performed experiments and shows the results obtained; Section 5 analyses and discusses the results; Section 6 draws conclusions and future work; finally, Appendixes A to C show abbreviations, deep learning architectures and tables that reinforce the results.

II. BACKGROUND

Section 2 presents the background and describes some of the works that relate FS to deep learning and multi-objective metaheuristics, as well as work related to surrogate-assisted evolutionary algorithms and air quality time series forecasting.

A. Feature selection

FS is the process of removing irrelevant and redundant attributes from a database [7]. The goal of FS is to find the best subset of attributes, reducing the cost and computational time required for model learning. Another interesting property is that it can improve the interpretability of the models by making them less complex and therefore easier to understand and interpret. A FS process can be either an attribute subset or an attribute evaluation process. The former searches for subsets of attributes using some search strategy. Their advantage is that they consider the interaction between attributes (multivariate methods), but the search requires a number of steps that is $O(2^n)$ (NP-hard problem) with n being the number of attributes. The latter evaluates attributes in order to assign them a ranking or importance, thus they are called *Feature Ranking* (FR) methods [8]. These methods are usually univariate, although there are a few multivariate methods, such as *relieff*

[9]. FR techniques can also be treated as FS techniques if a subset with the best attributes in the ranking, or those above a certain threshold of importance, is selected. As shown in [10], where Zhang *et al.* propose a $l_{0,2}$ -norm based FS method to select the top k features in various scenarios. The subset evaluation FS process has several stages. First, from the initial database, a subset of attributes is selected with a search strategy. Some of the most popular search strategies are *forward feature selection* [11], *recursive feature elimination* [12] or *multi-objective evolutionary search* [13]. The subsets generated by the search strategy are evaluated with some evaluation function. If after evaluation a stop criterion is satisfied, the process is terminated and moves to a validation phase. If the stop criterion is not satisfied, a new subset is generated again until the stop criterion is satisfied.

FS methods can typically be classified into three types: filter, wrapper and embedded. Filter methods evaluate attributes through statistical functions such as *correlation* [14], *consistency* [15], *redundancy* [16], *information gain* [17], χ^2 [18], etc. They are the fastest and simplest. Wrapper methods build machine learning models whose performance is evaluated using some metric, such as *accuracy* and *area under the receiver operating characteristic (ROC) curve* for classification tasks, or RMSE and *mean absolute error* (MAE) for regression. Therefore, the computational cost is usually high, although they are, in general, more accurate than filter FS methods. *Embedded* methods perform FS within the model training process. Examples of this can be *least absolute shrinkage and selection operator* (LASSO) [19], *C4.5* [20] or *classification and regression trees* (CART) [21]. *Ensembles* can also be used as an FS technique [22]. They are a combination of the results obtained by multiple techniques. Ensembles try to achieve better outcomes than would be expected using only one method. Finally, there are *hybrid methods* that generally combine filter and wrapper. The purpose of this is usually to improve the prediction error of the filters and reduce the computational cost of the wrapper.

FS in deep learning models is less common than in machine learning. However, in the literature, we can find several articles on FS with deep learning. Embedded methods tend to be the most common when performing FS with deep learning, as a layer can be added to the neural network to perform this process. For example, Borisov *et al.* [23] have developed a new input layer for *artificial neural networks* (ANN), called *CancelOut*, for FS in deep learning models. Other examples of neural networks that integrate FS are *NeuralFS* [24], based on DNN, and CNN-FS [25], based on *convolutional neural networks* (CNN) [26]. Wang *et al.* [27] combines CNN for feature extraction with χ^2 for FS, along with ensemble learning based on *extreme machine learning* (EML) and weighted voting. Attention-based methods are another effective method for handling features that have different importance in deep learning. Yuan *et al.* [28] propose a spatio-temporal attention-based LSTM network which can identify important input variables at each time step and adaptively discover hidden quality-related states at all time steps. Wang *et al.* [29] propose an integrated FS scheme based on a neural network with redundancy control and Group Lasso Regularization to

produce sparsity.

There are also modifications of embedded FS methods to consider discarded features such as the work of Zheng *et al.* [30], who propose an embedded framework for FS where an unselected feature classifier to find an optimal subset of attributes is added. In this line of research, there are other methods to retain more statistical and structural information about the features such as the case of Wu *et al.* [31] where an orthogonal least square regression model with feature weighting for FS is proposed.

FS has been successfully applied in environments where features are not previously known to the learner, but are streams whose information is increasing over time, as is the case in [32], where Zhou *et al.* propose a streaming feature selection considering feature interaction.

B. Multi-objective meta-heuristics for feature selection based deep learning

For the design of FS methods, the use of MOEAs is very popular. MOEAs simultaneously optimize several conflicting objective functions. This gives rise to the set of Pareto optimal solutions, i.e. solutions where one objective cannot be further improved without making another objective worse. The optimal solutions in a multi-objective optimization problem are called *non-dominated solutions*, and they form the *Pareto front*. Evolutionary algorithms are well suited for multi-objective optimization because they allow multiple non-dominated solutions to be captured in a single run of the algorithm, in addition to their ability to solve complex problems. MOEAs are called *many-objective evolutionary algorithms* [33] when solving optimization problems of more than three objectives.

Some works use multi-objective optimization algorithms based deep learning to perform FS. For example, the one proposed by Al-Tashi *et al.* [34] who present a wrapper based on ANN for FS in classification, with the *multi-objective grey wolf optimizer* (BMOGW) as the search strategy. Within the hybrid methods with deep learning, we can find combinations of filters and wrappers. For financial prediction, Niu *et al.* [35] combines *reliefF* filter and a wrapper based on EML, and BMOGW combined with *cuckoo search* to perform FS. Although multi-objective optimization is the most common, there are also studies focusing on many-objective optimization. Recently, Shu *et al.* [36] have proposed a 5-objective optimization model with *non-dominated sorting genetic algorithm III* (NSGA-III) [37], of which four are filters and the remaining one is a wrapper based on EML.

C. Surrogate-assisted evolutionary algorithms

Surrogate-assisted evolutionary computation attempt to approximate the fitness function in evolutionary algorithms through a more computationally efficient model [38]. This technique is especially useful when dealing with high-dimensional problems and has proven to be effective in multi-objective optimization problems [39], [40].

An example of a surrogate-assisted method with deep learning and multi-objective meta-heuristics used to perform

FS is shown in the work of Jiang *et al.* [41]. They propose an ensemble method based on a two-layer surrogate-assisted mechanism called *multi-surrogate-assisted dual-layer ensemble feature selection* (MDEFS). To select the most relevant features in large datasets, a subset with the most representative samples is selected with three strategies based on surrogate-assisted models: *k-means*, *density-based spatial clustering of applications with noise* and *random sub-sampling*. Then, a *particle swarm*-based algorithm is used to select the best attributes.

D. Air quality

Air quality forecasting with multi-objective evolutionary computation combined with deep learning is a field that is still under explored. However, there are already some papers that address this issue. Most of them are focused on the predictions of $PM_{2.5}$ in different areas of China, such as Liu *et al.* [42], that introduce a multi-resolution ensemble model based on *wavelet packet decomposition*, *bidirectional LSTM* and *nondominated sorting genetic algorithm II* (NSGA-II) [43], which tries to minimize bias and variance, to obtain deterministic forecasts. Zhang *et al.* [44] propose an ensemble model based on *multi-objective ensemble pruning with multi-population NSGA-II* and MLP, CNN and LSTM for stable $PM_{2.5}$ time series forecasting. Wang *et al.* [45] developed a hybrid system based on two-step FS with correlation, *reliefF* and a modified *quantum fuzzy neural network* for air quality forecasting. The relevance of the selected attributes is determined through a multi-objective chaotic map-based algorithm called *multi-objective chaotic bonobo optimizer*. Du *et al.* [46] propose a hybrid model, in this case to predict $PM_{2.5}$ and PM_{10} concentrations. Forecasts are made using *multi-objective Harris hawks optimization* and EML. The two objectives to optimize are prediction accuracy and stability of forecasting errors.

E. Conclusion of the related works

The study of the state of the art reveals that, although currently the subject of FS for deep learning is being of interest to researchers, it is necessary to deepen in certain aspects. For example, few studies have been carried out on the application of multivariate methods of filter-type FS for deep learning, both for search strategies and evaluators. Regarding wrapper methods for deep learning, although some studies have proposed approaches to reduce the computational cost, this remains a challenge for researchers. Finally, the embedded FS methods for deep learning have been compared fundamentally with other embedded methods such as LASSO or decision trees, but not so much with other powerful multivariate FS methods of filter, wrapper, hybrid or surrogate-assisted type. Therefore, in this paper we propose and discuss different FS multivariate hybrid methods for deep learning that combine filter methods with wrapper surrogate-assisted methods, using multi-objective evolutionary computation as a search strategy. The proposed methods are compared with a wide range of FS methods that include univariate, multivariate, filter, wrapper and embedded methods.

The authors of this paper have extensive experience in the fields of FS, multi-objective evolutionary computation, time series forecasting and air quality prediction. In [13] we proposed a MOEA, called ENORA, for FS in regression problems. ENORA was used as a search strategy in a wrapper method based on random forest. It was necessary to reduce the number of trees of the random forest algorithm during the optimization process to reduce the computational time of the method. In [47] and [48], MOEAs were used for FS with time series data applied to the prediction of antibiotic resistance outbreaks and the forecast of energy consumption in smart buildings, respectively. In these works, FS wrapper methods based on deep learning were prohibitive due to the excessive computational time required. In [49] deep learning techniques, particularly LSTM, GRU and CNN, are used for air quality forecasting with time series data. However, in this work FS was not performed, being proposed for future work. Finally, in [50] a MOEA was proposed for the spatio-temporal forecast of air pollution. The proposed technique builds a model based on ensemble learning with the multiple linear regression models found with the MOEA, which is compared with quasi-recurrent neural networks, among other models. Again, FS was raised for future work. Therefore, this paper is proposed as a continuation of the previous research, avoiding the inconvenience of computational costs required by deep learning-based FS techniques, by using surrogate-assisted MOEAs.

III. MATERIALS AND METHODS

In this section, the FS problem is formalized as a multi-objective Boolean optimization problem. The proposed multi-objective optimization models are mathematically defined and the main components of the surrogate-assisted evolutionary algorithm that solve them are described.

A. Feature selection as a multi-objective boolean optimization problem

FS can be formulated as a *multi-objective boolean optimization problem* as follows:

$$\text{Min./Max. } f_k(\mathbf{x}), \quad k = 1, \dots, l \quad (1)$$

where $f_k(\mathbf{x})$ are *objective functions* (wrapper type or filter type), $\mathbf{x} = \{x_1, x_2, \dots, x_w\} \in \mathbb{B}^w$ represents the set of *decision variables*, and $\mathbb{B} = \{\text{true}, \text{false}\}$ is the domain for each variable x_i , $i = 1, \dots, w$. In problem (1), $x_i = \text{true}$ represents that attribute i is selected, and $x_i = \text{false}$ represents that attribute i is not selected, for all $i = 1, \dots, w$. Let $X = \{\mathbf{x} \in \mathbb{B}^w\}$ be the *search space* of the problem (1). We want to find a subset of solutions $\Psi \subseteq X$ called *non-dominated set* (or *Pareto optimal set*). A solution $\mathbf{x} \in X$ is *non-dominated* if there is no other solution $\mathbf{x}' \in X$ that dominates \mathbf{x} , and a solution \mathbf{x}' *dominates* \mathbf{x} if and only if there exists k ($1 \leq k \leq l$) such that $f_k(\mathbf{x}')$ improves $f_k(\mathbf{x})$, and for every k ($1 \leq k \leq l$), $f_k(\mathbf{x})$ does not improve $f_k(\mathbf{x}')$. In other words, \mathbf{x}' *dominates* \mathbf{x} if and only if \mathbf{x}' is better than \mathbf{x} for at least one objective, and not worse than \mathbf{x} for

any other objective. For minimization problems, the set Ψ of non-dominated solutions of (1) can be formally defined as:

$$\Psi = \{\mathbf{x} \in X \mid \nexists \mathbf{x}' \in X \mid \mathcal{D}(\mathbf{x}', \mathbf{x})\}$$

where $\mathcal{D}(\mathbf{x}', \mathbf{x})$ is equivalent to:

$$\exists k, 1 \leq k \leq l, f_k(\mathbf{x}') < f_k(\mathbf{x}) \wedge \forall k, 1 \leq k \leq l, f_i(\mathbf{x}') \leq f_i(\mathbf{x})$$

Once the set of optimal solutions is available, the most satisfactory one can be chosen by applying a preference criterion. Search space $X = \{\mathbf{x} \in \mathbb{B}^w\}$ of problem (1) contains 2^w solutions (*NP-hard* problem). Metaheuristics methods such as evolutionary algorithms are typically used to find approximate solutions for NP-hard class problems, including FS problems [51], and they are particularly appropriate for multi-objective optimization.

B. Proposed multi-objective optimization models for feature selection

Let $D = \{\mathbf{d}_1, \dots, \mathbf{d}_s\}$ be a normalized dataset with s instances. Each instance $\mathbf{d}_t = \{d_t^1, \dots, d_t^w, o_t\}$, $t = 1, \dots, s$, has w input attributes of any type, and one output attribute $o_t \in \mathbb{R}$. In this approach, the initial dataset D is divided into three partitions R , V and T for training, validation and test with 60%, 20% and 20% of the data respectively. Then these three partitions are normalized. We consider multi-objective boolean optimization models for FS defined as in (1). Let $\mathcal{F}_{R,V}^\Phi(\mathbf{x})$ be a performance measure, e.g. RMSE, MAE, *correlation coefficient* (CC), etc., of a regression model built with a learning algorithm Φ trained with dataset $R \subset D$ and evaluated with validation dataset $V \subset D$ using only the selected attributes $x_i = \text{true}$, $i = 1, \dots, w$. Let $\mathcal{C}(\mathbf{x})$ be a function that measures the number of selected attributes of \mathbf{x} , i.e.:

$$\mathcal{C}(\mathbf{x}) = \sum_{i=1}^w \mathcal{N}(x_i) \quad (2)$$

where \mathcal{N} is a function that transforms a boolean value into numeric ($\text{true} = 1$ and $\text{false} = 0$).

Let $\mathcal{P}_R(\mathbf{x})$ be a function that measures the sum of correlation between each selected attribute $x_i = \text{true}$, $i = 1, \dots, w$, and the output attribute in dataset R , defined as follows:

$$\mathcal{P}_R(\mathbf{x}) = \sum_{\substack{i=1 \\ x_i=\text{true}}}^w \rho_i^R \quad (3)$$

where ρ_i^R is the normalized Pearson's correlation coefficient between the selected attribute i and the output in dataset R .

Let $\mathcal{S}_R(\mathbf{x})$ be a function that calculates the sum of *reliefF* scores of the selected attributes of \mathbf{x} in dataset R , i.e.:

$$\mathcal{S}_R(\mathbf{x}) = \sum_{\substack{i=1 \\ x_i=\text{true}}}^w \sigma_i^R \quad (4)$$

where σ_i^R is normalized *reliefF* score of attribute i in dataset R . In view of this, we consider the next four objectives:

Objective O1. Minimize the RMSE of the regression model built with learning algorithm Φ trained with dataset R and evaluated with validation dataset V :

$$\text{Minimize } O1 \equiv \mathcal{F}_{R,V}^{\Phi}(\mathbf{x})$$

Objective O2. Minimize the number of selected attributes:

$$\text{Minimize } O2 \equiv \mathcal{C}(\mathbf{x})$$

Objective O3. Maximize correlation between the selected attributes and the output variable in dataset R :

$$\text{Maximize } O3 \equiv \mathcal{P}_R(\mathbf{x})$$

Objective O4. Maximize sum of *reliefF* scores of the selected attributes in dataset R :

$$\text{Maximize } O4 \equiv \mathcal{S}_R(\mathbf{x})$$

The following *multi-objective boolean optimization problems* for modelling wrapper-filter FS methods in regression tasks are instances of problem (1) and have been proposed to optimize the objectives defined previously:

- *2-objective optimization model for wrapper FS method.* Objectives to optimize are $O1$ and $O2$. This problem will be called $O1O2$.

$$\begin{aligned} \text{Minimize } O1 &\equiv \mathcal{F}_{R,V}^{\Phi}(\mathbf{x}) \\ \text{Minimize } O2 &\equiv \mathcal{C}(\mathbf{x}) \end{aligned} \quad (5)$$

- *3-objective optimization model for filter FS method based on correlation and reliefF.* Objectives to optimize are $O2$, $O3$ and $O4$. This problem will be called $O2O3O4$.

$$\begin{aligned} \text{Minimize } O2 &\equiv \mathcal{C}(\mathbf{x}) \\ \text{Maximize } O3 &\equiv \mathcal{P}_R(\mathbf{x}) \\ \text{Maximize } O4 &\equiv \mathcal{S}_R(\mathbf{x}) \end{aligned} \quad (6)$$

- *3-objective optimization model for wrapper-filter FS method based on correlation.* Objectives to optimize are $O1$, $O2$ and $O3$. This problem will be called $O1O2O3$.

$$\begin{aligned} \text{Minimize } O1 &\equiv \mathcal{F}_{R,V}^{\Phi}(\mathbf{x}) \\ \text{Minimize } O2 &\equiv \mathcal{C}(\mathbf{x}) \\ \text{Maximize } O3 &\equiv \mathcal{P}_R(\mathbf{x}) \end{aligned} \quad (7)$$

- *3-objective optimization model for wrapper-filter FS method based on reliefF.* Objectives to optimize are $O1$, $O2$ and $O4$. This problem will be called $O1O2O4$.

$$\begin{aligned} \text{Minimize } O1 &\equiv \mathcal{F}_{R,V}^{\Phi}(\mathbf{x}) \\ \text{Minimize } O2 &\equiv \mathcal{C}(\mathbf{x}) \\ \text{Maximize } O4 &\equiv \mathcal{S}_R(\mathbf{x}) \end{aligned} \quad (8)$$

- *4-objective optimization model for wrapper-filter FS method based on correlation and reliefF.* Objectives to optimize are $O1$, $O2$, $O3$ and $O4$. This problem will be called $O1O2O3O4$.

$$\begin{aligned} \text{Minimize } O1 &\equiv \mathcal{F}_{R,V}^{\Phi}(\mathbf{x}) \\ \text{Minimize } O2 &\equiv \mathcal{C}(\mathbf{x}) \\ \text{Maximize } O3 &\equiv \mathcal{P}_R(\mathbf{x}) \\ \text{Maximize } O4 &\equiv \mathcal{S}_R(\mathbf{x}) \end{aligned} \quad (9)$$

C. Surrogate-assisted multi-objective evolutionary algorithm

The function $\mathcal{F}_{R,V}^{\Phi}(\mathbf{x})$ measures the performance (we use RMSE) of a prediction model training with learning algorithm Φ (we use LSTM) on dataset R with attributes of \mathbf{x} and evaluated on validation dataset V . $\mathcal{F}_{R,V}^{\Phi}(\mathbf{x})$ can be used as evaluator in a wrapper FS method. As is well known, the drawback of wrapper methods for FS is the computational time required by the learning algorithm to build the models with the attribute subsets explored with the search strategy. This disadvantage is increased in the presence of large datasets or when the learning algorithm is particularly expensive, as is the case with deep learning. To take advantage of the good properties of wrapper methods without the need to train the model on each subset of candidate attributes, in this paper we propose a *surrogate-assisted* approach. The set of attributes selected by the surrogate-assisted MOEA proposed in this paper is finally evaluated on the test set T , which has not been seen by the MOEA. Figure 1 shows the flowchart of the proposed surrogate-assisted MOEA, and its main components are shown below.

1) *Representation of solutions:* Our MOEA uses a fixed-length representation, where each individual consists of a bit vector for attribute selection. Therefore, an individual I is represented as:

$$I = \{b_1^I, \dots, b_w^I\}$$

where $b_i^I \in \{0, 1\}$ for $i = 1, \dots, w$. Each bit $b_i^I \in \{0, 1\}$ represents an attribute in the dataset (1 for selected, and 0 for non-selected attributes).

2) *Initial population:* The initial population is randomly generated with a Bernoulli distribution with probability 0.5.

3) *Fitness function:* We use the surrogate function $\mathcal{Q}_{R,V}^{\Phi}(\mathbf{x})$ as an approximation to the function $\mathcal{F}_{R,V}^{\Phi}(\mathbf{x})$ for the calculation of the objective function $O1$. $\mathcal{Q}_{R,V}^{\Phi}(\mathbf{x})$ is the RMSE of a surrogate prediction model M_R^{Φ} trained with learning algorithm Φ (we use LSTM) and dataset R considering all w attributes, and evaluated with a dataset V' which consists of the validation dataset V where attribute values $x_i = false$, $i = 1, \dots, w$, are set to a constant value α in all s instances (we used $\alpha = 0$ in the experiments, although other values can also be used). Note that surrogate model M_R^{Φ} does not depend on the set \mathbf{x} of decision variables, so it is trained only once in off-line mode. Note also that $\mathcal{Q}_{R,V}^{\Phi}(\mathbf{x}_D) = \mathcal{F}_{R,V}^{\Phi}(\mathbf{x}_D)$, where \mathbf{x}_D is the solution vector with $x_i = true$, for all $i = 1, \dots, w$. Algorithm 1 shows the calculation of the fitness function for objective $O1$. The fitness functions for the rest of the objective functions $O2$, $O3$ and $O4$ correspond to the mathematical formulations described in equations (2), (3) and (4) respectively.

4) *Variation operators:* The variation operators used are *half uniform crossover* [52] and *bit flip mutation* [53]. We have chosen these operators since they are the ones used by the Platypus platform (<https://platypus.readthedocs.io/en/latest/>) for MOEAs with binary representation.

5) *Decision process to choose the final non-dominated solution:* For each non-dominated individual in the last population, an LSTM model m is trained on dataset R with the attributes selected in the individual, and a *recursive multi-step strategy*

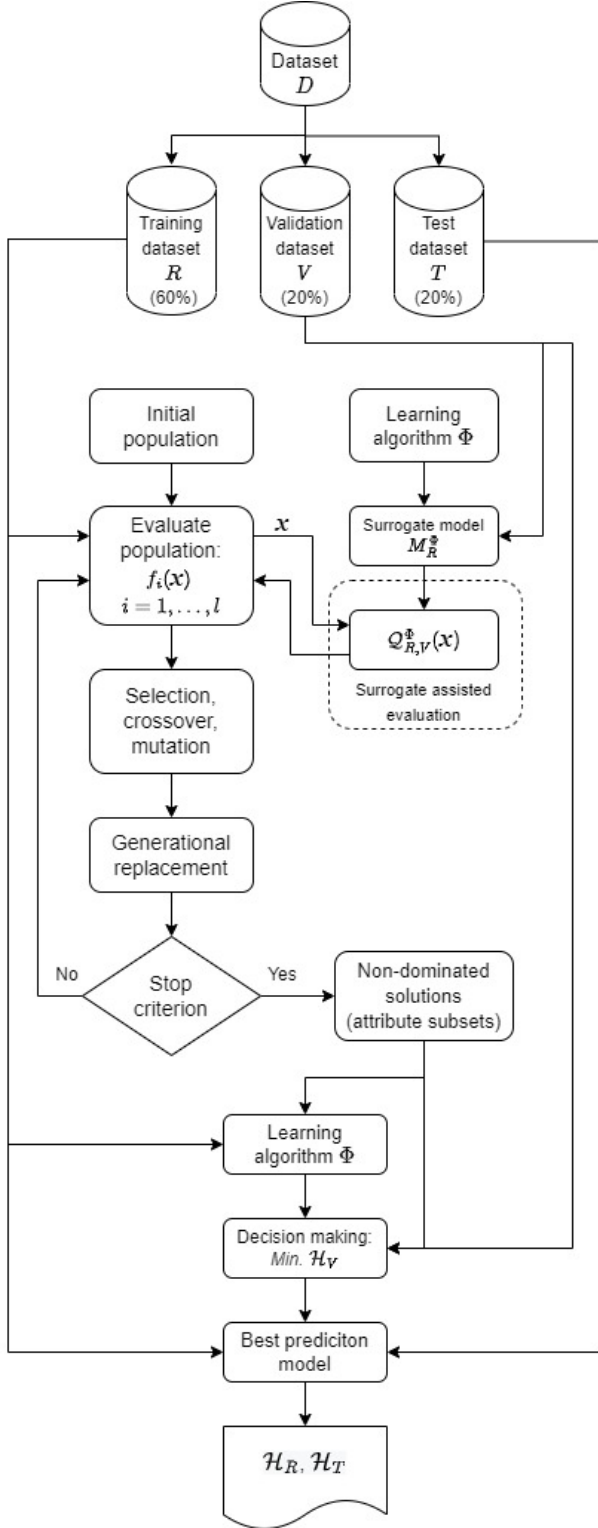


Fig. 1: Flowchart of the proposed surrogate-assisted multi-objective evolutionary algorithm.

Algorithm 1 Fitness function for objective function OI

Require: $I = \{b_1^I, \dots, b_w^I\}$ {Individual}

Require: M_R^Φ {Surrogate prediction model built with learning algorithm Φ and trained with dataset R with all attributes}

Require: $V \subset D$ {Validation dataset}

Require: α {Imputation constant}

1: $V' \leftarrow V$

2: **for** $i = 1$ **to** w **do**

3: **if** $b_i^I = 0$ **then**

4: **for** $d_t \in V'$ **do**

5: $d_t^i = \alpha$

6: **end for**

7: **end if**

8: **end for**

9: **return** $RMSE(M_R^\Phi, V')$ {RMSE of surrogate prediction model M_R^Φ evaluated in dataset V' }

[54] is performed for h -steps ahead. Although there are other strategies for multi-step ahead time series forecasting [55], [56], in this research we have decided to use the recursive strategy. The recursive strategy constrains all the horizons to be forecasted with the same model structure, and can suffer from poor performance in multi-step forecasting tasks, especially if the forecast horizon h exceeds the window size l , since at some point all the inputs are forecast values rather than actual observations. Nevertheless, this does not occur in our research since we use $h \leq l$. In spite of these drawbacks, we have preferred to adopt the recursive strategy as a first approximation of our work since it is not computationally expensive and can be implemented in all the forecasting methods compared in this paper.

The non-dominated individual are evaluated on dataset V using a performance measure \mathcal{H} . The chosen non-dominated individual will be the one with the lowest value of \mathcal{H} . Therefore, an ideal solution is one with $\mathcal{H} = 0$. Algorithm 2 calculates the performance measure \mathcal{H} of a solution x on a test dataset E , which is particularly used to evaluate of the non-dominated individuals for MOEA decision making by setting $E = V$.

IV. EXPERIMENTS AND RESULTS

This section describes the air quality dataset used in the experiments and its preprocessing. The experiments have been established for the comparison of the different multi-objective optimization models for FS proposed in this paper, the comparison of the MOEAs that solve them, the identification and adjustment of the best prediction model and the comparison with other FS techniques.

A. Air quality dataset

Data are extracted from the *Autonomous Community of the Region of Murcia, Spain* (<https://sinqlair.carm.es/calidadaire/redvigilancia/redvigilancia.aspx>), which provides information on air quality in the Region of Murcia thanks to *Consejería de Agua, Agricultura, Ganadería, Pesca y Medio Ambiente*.

Algorithm 2 Evaluation \mathcal{H} of a solution on a dataset E

Require: $x = \{x_1, \dots, x_w\}$ {Solution}
Require: m {Prediction model trained with the attributes selected in the solution x }
Require: E {Dataset for evaluation}
Require: h {Number of steps-ahead}
Require: $w_1, w_2, w_3, 0 \leq w_1, w_2, w_3 \leq 1, w_1 + w_2 + w_3 = 1$ {Weights}

- 1: $E' \leftarrow E$
- 2: **for** $i = 1$ **to** w **do**
- 3: **if** $x_i = 0$ **then**
- 4: $\text{remove}(E', i)$ {Remove attribute i from dataset E' }
- 5: **end if**
- 6: **end for**
- 7: $RMSE \leftarrow \frac{1}{h} \sum_{j=1}^h RMSE(m, j, E')$
- 8: $MAE \leftarrow \frac{1}{h} \sum_{j=1}^h MAE(m, j, E')$
- 9: $CC \leftarrow \frac{1}{h} \sum_{j=1}^h |1 - CC(m, j, E')|$
- 10: $\mathcal{H} \leftarrow w_1 \cdot RMSE + w_2 \cdot MAE + w_3 \cdot CC$
- 11: **return** \mathcal{H} {Evaluation of a solution}



Fig. 2: Photograph of the air quality measurement station in La Aljorra. Google image.

All the collected data comes from La Aljorra station (Figure 2), taken daily for four years, from 2017 to 2020. In total, there are 1461 instances. The initial dataset consisted of 17 columns: $Date$, NO , NO_2 , SO_2 , O_3 , TMP (temperature), HR (relative humidity), NO_X , DD (wind direction), PRB (atmospheric pressure), RS (solar radiation), VV (wind speed), C_6H_6 , C_7H_8 , XIL , PM_{10} , $Noise$. Table I shows a summary of the values of all attributes.

B. Data preprocessing: missing values imputation and sliding window transformation

The proposed dataset is used for NO_2 prediction. For the initial preprocessing of the data, those columns with more than 25% missing values were removed. These columns were: C_6H_6 , C_7H_8 , XIL and $Noise$. $Date$ column has also been eliminated, as it does not provide any relevant information to

TABLE I: Summary of initial attributes.

Name	Units	Count	Mean	Std	Min	Max
$Date$	-	1461	-	-	-	-
NO	$\mu g/m^3 N$	1272	4.38	2.43	1.00	31.00
SO_2	$\mu g/m^3 N$	1299	9.21	3.45	2.00	23.00
O_3	$\mu g/m^3 N$	1403	57.94	15.25	19.00	112.00
TMP	$^{\circ}C$	1409	19.59	5.74	4.00	32.00
HR	% R.H.	1409	69.21	14.09	27.00	100.00
NO_X	$\mu g/m^3 N$	1272	21.21	11.20	3.00	104.00
DD	degrees	1409	192.11	104.45	0.00	360.00
PRB	mb	1409	1017.81	6.41	991.00	1033.00
RS	W/m^3	1409	182.32	82.74	13.00	338.00
VV	m/s	1409	1.22	0.43	1.00	3.00
C_6H_6	$\mu g/m^3 N$	410	0.65	0.51	0.10	3.40
C_7H_8	$\mu g/m^3 N$	410	0.93	0.40	0.30	3.10
XIL	$\mu g/m^3 N$	410	1.40	0.68	0.10	3.00
PM_{10}	$\mu g/m^3 N$	1381	26.27	12.98	5.00	168.00
$Noise$	dba	0	-	-	-	-
NO_2 (target)	$\mu g/m^3 N$	1272	14.64	8.19	2.00	58.00

the treated problem. The remaining missing values have been imputed with linear interpolation.

Lagged transformation of the input variables has to be done with *sliding window* method [57] in order to remove time dependencies in the data. Let l be the *window size* (WS), i.e., the number of previous time steps. The sliding window transformation process builds the following dataset D_l :

$$D_l = \{ \{d_{t-1}^1, \dots, d_{t-l}^1\}, \dots, \{d_{t-1}^w, \dots, d_{t-l}^w\}, \{o_{t-1}, \dots, o_{t-l}\}, o_t \},$$

$$t = 1, \dots, s \quad (10)$$

Note that $d_t^i, i = 1, \dots, w$, and o_t values with $t \leq 0$ do not exist and are therefore considered missing values. The new attributes are named $\langle Lag \rangle_ \langle name \text{ of the original attribute} \rangle_ \langle number \text{ of previous time step} \rangle$. For example, NO attribute with 1 previous time step will be called Lag_NO_1 . A window size of 7 has been selected, representing one week. The rows with missing values resulting from the sliding window transformation have been removed. Note that the original variables $d_t^i, i = 1, \dots, w$, of the input attributes have been removed, as it is not correct to use the future value of an attribute to obtain future predictions. The only original variable held is the output. Therefore, our dataset with WS 7 has 1454 rows and 85 columns. Finally, the dataset has been split into 60% for training, 20% for validation and 20% for testing, preserving the order of the instances and then normalized.

C. Multi-objective optimization models comparison

This section describes the process of searching for the best hyper-parameters of the surrogate deep learning model, the evaluation of the prediction models for each multi-objective optimization model and the statistical tests applied to select the best model.

1) *Hyper-parameter tuning of the surrogate-assisted model:* For the creation of the M_R^Φ surrogate prediction model used in the evolutionary multi-objective algorithm we have selected LSTM. Our previous study [49] showed that an LSTM with one hidden layer with *rectified linear unit* [58] activation

function and a dense layer with a linear activation as the output layer can obtain good results in time series forecasting applied to air quality prediction, over other learning algorithms such as 1D-CNN, random forest, SVM or LASSO regression. In order to optimize this model, we first search for the best hyper-parameters. For this purpose, a 3-fold cross-validation with one repetition was used. A grid search has been carried out with: 1, 2 and 3 hidden layers; 50, 100, 150, 200, 250 and 300 neurons; a batch size of 32, 64, 128, 512 and 1024; and 100, 500 and 1,000 epochs. Therefore, 270 possible combinations have been performed. The best combination of hyper-parameters obtained an RMSE of 0.07877 with a standard deviation of 0.01339 and was 1 hidden layer, 200 neurons, batch size of 128 and 100 epochs. The LSTM architecture used as surrogate model in this paper consists of an input layer, a hidden layer with 200 neurons and relu activation function, a layer with a dropout of 0.2 and an output layer with one neuron and linear activation function.

A surrogate LSTM model M_R^Φ has been trained with all attributes, a 60% of the instances and with the best hyper-parameters, and tested with 20% of the validation instances. The RMSE in training dataset is 0.0515 and the RMSE in validation dataset is 0.1795.

2) *Prediction model evaluation for each multi-objective optimization model:* For our dataset, the proposed multi-objective optimization problems (5), (6), (7), (8) and (9) have been solved with NSGA-II. NSGA-II is widely recognized by the scientific community as representative algorithm for solving multi-objective optimization problems. This MOEA is used in this initial phase of the experiments to compare the different multi-objective optimization models, although in a later experimentation phase (section IV-D) NSGA-II is compared with a broader set of MOEAs. It has been run 10 times for each problem, with population size of 100, 1,000 generations (100,000 evaluations of the objective functions) and crossover and mutation probabilities of 1.0. Experiments were run on a computer with an AMD Ryzen 7 5800X 8-Core Processor with 3.80 GHz using 16 GB of RAM and Windows 10 Pro. We use RMSE as performance metric $Q_{R,V}^\Phi(x)$. For the decision making process in order to choose the best non-dominated solution (Algorithm 2 with the validation dataset V) we have used multi-step ahead predictions and normalized weights $w_1 = w_2 = w_3$. Finally, the best model obtained with NSGA-II is evaluated by using Algorithm 2 with the (unseen) test dataset T .

Table II shows for each optimization problem, the results of the best solution including the evaluation \mathcal{H} on training, validation and test datasets (\mathcal{H}_R , \mathcal{H}_V and \mathcal{H}_T respectively) and the number of selected attributes. Therefore, we will have 6 datasets: 5 reduced datasets selected by the NSGA-II algorithm and the original dataset.

3) *Statistical tests for performance prediction models:* Once the best reduced dataset for each optimization problem has been found, 3 repetitions of 10-fold cross-validation with the training dataset R (a total of 30 prediction models) are performed using the Algorithm 2 for evaluation. In order to apply statistical tests, first, we check whether the conditions of normality and sphericity of the samples are met. For normality,

TABLE II: Prediction model evaluation with NSGA-II for the air quality problem, 100,000 evaluations and 10 runs, sorted from best to worse evaluation of \mathcal{H}_T .

Optimization model	$Q_{R,V}^\Phi(x)$	Number of selected attributes	\mathcal{H}_R	\mathcal{H}_V	\mathcal{H}_T	Run time (minutes)
<i>O1O2O3</i>	0.1939	16	0.0807	0.2182	0.1298	15.76
<i>O1O2</i>	0.2403	2	0.1234	0.2328	0.1442	9.80
<i>O1O2O3O4</i>	0.2006	4	0.1210	0.2347	0.1447	28.41
<i>O1O2O4</i>	0.2462	6	0.0928	0.2447	0.1647	18.33
<i>O2O3O4</i>	–	19	0.1006	0.2746	0.1965	7.65

Shapiro-Wilks test was used, if p -values are greater than 0.05 then it can be assumed that the values come from a normal distribution. For sphericity we used Mauchly's test, if p -values are greater than 0.05 then the sphericity condition is satisfied. Then, tests are applied to check if there are statistically significant differences between the models. Parametric ANOVA test in case the normality and sphericity conditions are met and non-parametric Friedman test in case they are not met. If Friedman test is applied, Nemenyi post hoc test is also used to determine where the significant differences between the models are (those with a p -value less than 0.05).

In order to summarize results and determine the best and worst models, method-dataset pairs have been ranked according to the difference between *wins* and *losses* obtained by Wilcoxon signed-rank test. Every time one method-dataset pair tests statistically significantly better than another, it counts as a *win* and otherwise as a *loss*. Table III shows the ranking of every multi-objective optimization model.

TABLE III: Ranking of multi-objective optimization models, 10 folds cross-validation, 3 repetitions, sorted from best to worse.

Optimization model	Win	Loss	Win-Loss
<i>O1O2O3</i>	4	0	4
<i>O1O2O4</i>	3	1	2
<i>O1O2</i>	1	2	-1
<i>O2O3O4</i>	1	2	-1
<i>O1O2O3O4</i>	0	4	-4

D. Multi-objective evolutionary algorithms comparison

1) *MOEAs performance evaluation:* Once the best optimization model has been identified, the best MOEA is determined in this new phase of experiments. We have selected the MOEAs implemented in the Platypus platform that allow binary representations, which are NSGA-II, NSGA-III, *multi-objective evolutionary algorithm based on decomposition* (MOEA/D) [59], *indicator-based evolutionary algorithm* (IBEA) [60], ϵ -MOEA [61], *strength pareto evolutionary algorithm 2* (SPEA2) [62] and ϵ -NSGA-II [63]. Each algorithm is run 30 times with the best multi-objective model (*O1O2O3*) with 100,000 evaluations and probability 1.0 of crossover and mutation. Afterwards, *hypervolume* metric is calculated. A set of *reference points* is required for the calculation of hypervolume metric. For the multi-objective optimization model *O1O2O3*, the set of reference points consists of 84 points (as many as attributes in the original dataset) in a three-dimensional space (one dimension for each objective). Dimension *O1* with 0 attributes is the value of $Q_{R,V}^\Phi(x_D)$

model when no attribute is selected, i.e. 0.2647 in our case. All other values of that dimension are 0. Dimension $O2$ indicates the number of attributes, so it will range from 0 to 84. Dimension $O3$ is the sum of the list ordered from best to worst correlation. That is, for 0 attributes it will be 0, for 1 attribute it will be the best correlation, for 2 attributes it will be the sum of the two best correlations and so on.

2) *Statistical tests for MOEA's performance:* The statistical test Mann-Whitney U rank is applied in order to rank the algorithms and select the best. Table IV shows the win-loss ranking for MOEAs.

TABLE IV: Ranking of the optimization algorithms with hypervolume metric, 100,000 evaluations, 30 runs, sorted from best to worse.

MOEA	Win	Loss	Win-Loss
NSGA-II	6	0	6
IBEA	4	1	3
ϵ -MOEA	4	1	3
ϵ -NSGA-II	3	3	0
NSGA-III	2	4	-2
SPEA2	1	5	-4
MOEA/D	0	6	-6

E. Final prediction model identification and adjustment

For each metric, those MOEAs that have never lost in any of the statistical tests are selected. The chosen algorithm is NSGA-II. Therefore, the best model is found with $O1O2O3$ optimization model with NSGA-II algorithm with crossover and mutation probabilities 1.0.

1) *Hyper-parameter tuning of the MOEA:* A hyper-parameter tuning is carried out on NSGA-II to establish the best crossover and mutation probabilities. The probabilities of both crossover and mutation on which the search has been performed have been 0.2, 0.5, 0.8 and 1.0 with 10 runs and 10,000 evaluations. The best combination was a crossover probability of 0.2 and a mutation probability of 0.2 for the hypervolume metric.

2) *Model adjustment:* The best MOEA (NSGA-II) is rerun for the best optimization model ($O1O2O3$) with the seed that produced best result in the previous phase of experiments, and the best probabilities obtained for the hypervolume metric (crossover probability 0.2 and mutation probability 0.2) with 1,000,000 evaluations. From the last population, the non-dominated solutions are obtained, the models are trained with LSTM and the average of the 7-steps ahead predictions are calculated. The evaluation \mathcal{H} on the training dataset R , the validation dataset V and the test dataset T are $\mathcal{H}_R = 0.1234$, $\mathcal{H}_V = 0.2328$ and $\mathcal{H}_T = 0.1441$, respectively. Note that the prediction model obtained using crossover probability 0.2 and mutation probability 0.2 does not improve the prediction model obtained using crossover and mutation probability 1.0. This is because the hypervolume metric measures the optimality and diversity of the MOEA based on the ideal population, but this metric do not explicitly measure the performance of the individuals as defined in Algorithm 2. Therefore, the best model with probability 1.0 of crossover and mutation is considered.

F. Comparison with other FS techniques

In this section, we compare our proposed solution with other FS techniques of different types. Among them, there are two filter-wrapper hybrid FS methods that combine attribute evaluation methods with LSTM and deterministic search, based on correlation and the reliefF algorithm respectively. We also compare our method with two wrapper multi-objective evolutionary FS methods (without surrogates) based on linear regression and random forest respectively, which solve the $O1O2$ multi-objective optimization model using NSGA-II. Finally, we compare two embedded FS methods: CancelOut, implemented in an LSTM, and random forest.

M1) Hybrid filter-wrapper FS method based on correlation and LSTM with deterministic search: This method, also proposed by the authors of this paper, uses the normalized Pearson's correlation coefficient between the selected attribute i and the output to sort all attributes from highest to lowest correlation. Then, w LSTM models are built with the i best attributes in the ranking, for all $i = 1, \dots, w$. The attribute subset that produces the LSTM model with the lowest RMSE is selected (Algorithm 3). This method is a hybrid between filter and wrapper since it uses together a correlation filter and an LSTM-based learning algorithm. It is also a multivariate FS method since although the attributes are ordered by their individual correlation with the output, attribute subsets are then formed with the most correlated attributes, in a deterministic way with a total of w candidate subsets, which are evaluated by means of an LSTM detecting, therefore, interactions between factors. Note that this method is equivalent to solving the multi-objective optimization problem formed by objectives $O2$ and $O3$, whose Pareto front consists of the set of solutions $\Psi = \{s_i \in \mathbb{B}^w, i = 1, \dots, w\}$ obtained with Algorithm 3, to finally select the solution s_i that builds the best prediction model m_{best} . However, the $O2O3$ multi-objective optimization problem can be solved deterministically in this case with Algorithm 3 without resorting to approximate probabilistic methods such as MOEAs.

M2) Hybrid filter-wrapper FS method based on reliefF and LSTM with deterministic search: This method is similar to the previous one except that it uses the multivariate reliefF filter instead of the univariate correlation filter, and has also been proposed by the authors of this paper. Algorithm 3 therefore also serves to describe this FS method. This method also uses a hybrid wrapper-filter technique, but in this case, it can be considered a fully multivariate method. This method is equivalent to solving the multi-objective optimization problem with the objectives $O2$ and $O4$ and selecting the best solution from the Pareto front, which can also be solved deterministically with Algorithm 3.

M3) Wrapper multi-objective evolutionary FS method based on linear regression: This is a multivariate wrapper method, for evaluating subsets of attributes, that uses NSGA-II with search strategy and a linear regression algorithm as evaluator together with the RMSE metric. The multi-objective optimization model $O1O2$ described in equation (5) is solved. This method does not use surrogates but rather a linear regression model is trained for each candidate attribute subset. It has been evaluated with 10 runs, 100,000 evaluations and 1.0 crossover

Algorithm 3 Hybrid filter-wrapper FS method based on correlation/relieff and LSTM with deterministic search

Require: $\mathbf{r} = \{r_1, \dots, r_w\}$, $1 \leq r_i \leq w \mid \rho_{r_{i-1}}^R \geq \rho_i^R, \forall i = 2, \dots, w$ $\{\sigma_{r_{i-1}}^R \geq \sigma_i^R, \forall i = 2, \dots, w$, for reliefF}

Require: $R \subset D$ {Training dataset}

Require: $V \subset D$ {Validation dataset}

Require: h {Number of steps-ahead}

Require: w_1, w_2, w_3 , $0 \leq w_1, w_2, w_3 \leq 1$, $w_1 + w_2 + w_3 = 1$ {Weights}

1: **for** $i = 1$ **to** w **do**

2: $x_i = 0$

3: **end for**

4: **for** $i = 1$ **to** w **do**

5: $x_{r_i} \leftarrow 1$

6: $m_i \leftarrow \text{LSTM}(R, \mathbf{x})$ {Prediction model built with LSTM in data set R with the attribute subset \mathbf{x} }

7: $s_i \leftarrow \mathbf{x}$

8: **end for**

9: $best \leftarrow \arg \min_{i=1, \dots, w} \mathcal{H}(s_i, m_i, V, h, w_1, w_2, w_3)$ {Using Algorithm 2}

10: **return** m_{best}

and mutation probabilities. Each set of selected attributes is evaluated with a linear regression model.

M4) Wrapper multi-objective evolutionary FS method based on random forest: This FS method is similar to the previous one except that the random forest learning algorithm is used for the evaluation of attribute subsets.

M5) CancelOut: CancelOut [23] is an embedded method that adds new input layer for ANNs. Each neuron in this layer is connected to an input variable. The weights are updated during a training stage so that irrelevant features are “cancelled” and those that do provide information are maintained. The CancelOut layer has three parameters: an activation function, λ_1 and λ_2 . λ_1 and λ_2 are two regularization terms to speed up the FR process. The values of both lambdas are between 0 and 1. A search has been performed to find the best CancelOut hyper-parameters. For this purpose, the CancelOut layer has been added to the previously described LSTM architecture. The values of λ_1 and λ_2 on which the search will be done are: 0, 0.0001, 0.0005, 0.001, 0.005, 0.01, 0.05, 0.1, 0.5, 1. The best value for λ_1 was 0.5 and for λ_2 was 0.05. The activation function will always be the sigmoid, which is the default parameter proposed by the authors of the layer. The LSTM architecture consists of an input layer, followed by a CancelOut layer, a hidden layer with 200 neurons and relu activation function, a layer with a dropout of 0.2 and an output layer with one neuron and linear activation function.

M6) Random Forest: Random forest is also an embedded method that performs feature selection internally. To do so, random forest considers how each attribute affects the impurity of the nodes.

In Table V, the compared FS techniques are sorted from best to worst evaluation of \mathcal{H} on the test dataset T (\mathcal{H}_T). The run times are in minutes. In addition, the original LSTM

model trained with all attributes has also been included in the comparison. Table VI shows the RMSE, MAE and CC of the 7-steps ahead predictions of the best prediction model (*OIO2O3-NSGA-II*) evaluated on the training dataset R and the test dataset T .

TABLE V: Comparison of feature selection methods for the airquality problem, sorted from best to worse evaluation of \mathcal{H}_T .

Method	\mathcal{H}_R	\mathcal{H}_T	Number of selected attributes	Run time (minutes)
<i>OIO2O3-NSGA-II</i>	0.0807	0.1298	16	15.76
<i>M1</i>	0.1246	0.1437	2	3.93
<i>M2</i>	0.1246	0.1437	2	4.09
<i>M4</i>	0.1235	0.1876	6	22.60
<i>M6</i>	0.1701	0.2069	1	2.16
<i>M3</i>	0.1227	0.2243	14	5.45
<i>M5</i>	0.0602	0.2452	84	0.07
<i>All attributes</i>	0.0560	0.2763	84	0.01

TABLE VI: Results of the best prediction model (obtained with *OIO2O3-NSGA-II* method) for the air quality problem evaluated on R and T datasets.

Set	Metric	1-step ahead	2-steps ahead	3-steps ahead	4-steps ahead	5-steps ahead	6-steps ahead	7-steps ahead
R	RMSE	0.0761	0.0744	0.0749	0.0751	0.0755	0.0758	0.0773
	MAE	0.0533	0.0529	0.0534	0.0535	0.0537	0.0538	0.0548
	CC	0.8892	0.8916	0.8900	0.8885	0.8861	0.8846	0.8805
T	RMSE	0.0814	0.0819	0.0822	0.0829	0.0832	0.0848	0.0873
	MAE	0.0494	0.0498	0.0503	0.0505	0.0507	0.0517	0.0537
	CC	0.7535	0.7508	0.7507	0.7478	0.7467	0.7380	0.7257

G. Additional experiments: indoor temperature forecasting in a domotic house

In this section we analyze a second forecasting problem in order to strengthen the conclusions about the proposal. The forecast of the indoor temperature in a domotic house in Valencia (Spain) is now considered. The dataset (*SML2010*) has been extracted from the *UCI Machine Learning Repository* [64]. In its original version [65], the dataset contains times series with a total of 24 attributes and 4137 instances, with data acquired in Valencia (Spain) from 03/13/2012 to 05/02/2012, each 15 minutes. In this study we have set the indoor temperature of the dining room as the target. After removing date and time attributes, useless attributes, and attributes related to the indoor temperature of another room, the set of 15 attributes shown in Tables VII and VIII is obtained.

Applying the methodology followed in this paper with a WS of 4, we obtain that the best optimization model has been *OIO2O3O4* in this case (Table IX). Note that the sliding window has not been applied to the DW variable, due to its categorical character. The best MOEA has turned out to be NSGA-II also in this problem. The *OIO2O3O4-NSGA-II* prediction model has been the best compared to the rest of the models obtained with the FS methods described in section IV-F. Table X shows the results of the comparison. Table XI shows the RMSE, MAE and CC of the 4-steps ahead predictions of the *OIO2O3O4-NSGA-II* model evaluated on the training dataset R and the test dataset T .

TABLE VII: Attributes in SML2010 dataset used in this paper.

Short name	Description
<i>WT</i>	Weather forecast temperature, in $^{\circ}C$
<i>CO₂</i>	Carbon dioxide, in <i>ppm</i> (dinning room)
<i>RH</i>	Relative humidity (dinning room), in %
<i>L</i>	Lighting (dinning room), in <i>Lux</i>
<i>R</i>	Rain, in range [0, 1]
<i>SD</i>	Sun dusk
<i>W</i>	Wind, in <i>m/s</i>
<i>SLW</i>	Sun light in west facade, in <i>Lux</i>
<i>SLE</i>	Sun light in east facade, in <i>Lux</i>
<i>SLS</i>	Sun light in south facade, in <i>Lux</i>
<i>SI</i>	Sun irradiance, in W/m^2
<i>OT</i>	Outdoor temperature, in $^{\circ}C$
<i>ORH</i>	Outdoor relative humidity, in %
<i>DW</i>	Day of the week, 1 = Monday, 7 = Sunday
<i>IT</i> (target)	Indoor temperature (dinning-room), in $^{\circ}C$ (target)

TABLE VIII: SML2010 attribute statistics.

Name	Mean	Std	Min	Max
<i>WT</i>	15.09	4.38	0.00	29.00
<i>CO₂</i>	206.60	22.76	187.34	594.39
<i>RH</i>	42.39	7.22	26.17	60.96
<i>L</i>	28.97	25.68	10.74	111.80
<i>R</i>	0.04	0.19	0.00	1.00
<i>SD</i>	335.09	304.51	0.61	625.00
<i>W</i>	1.30	1.22	0.00	6.32
<i>SLW</i>	14749.15	25306.45	0.00	95278.40
<i>SLE</i>	13566.28	23311.85	0.00	92367.50
<i>SLS</i>	19857.18	29494.60	0.00	95704.40
<i>SI</i>	232.20	312.46	-4.16	1094.66
<i>OT</i>	18.02	4.29	9.22	29.91
<i>ORH</i>	53.25	13.51	22.25	83.81
<i>DW</i>	3.96	1.99	1.00	7.00
<i>IT</i> (target)	20.49	3.31	11.35	28.92

TABLE IX: Prediction model evaluation with NSGA-II for the indoor temperature problem, 100,000 evaluations, 10 runs, sorted from best to worse evaluation of \mathcal{H}_T .

Optimization model	$\mathcal{Q}_{R,V}^{\Phi}(x)$	Number of selected attributes	\mathcal{H}_R	\mathcal{H}_V	\mathcal{H}_T	Run time (minutes)
<i>O1O2O3O4</i>	0.0364	10	0.0061	0.0091	0.0115	38.03
<i>O1O2O3</i>	0.0449	21	0.0075	0.0105	0.0130	27.18
<i>O1O2</i>	0.0659	4	0.0114	0.0126	0.0136	13.24
<i>O1O2O4</i>	0.0659	4	0.0114	0.0126	0.0136	32.34
<i>O2O3O4</i>	-	26	0.0062	0.0101	0.0149	14.43

TABLE X: Comparison of feature selection methods for the indoor temperature problem, sorted from best to worse evaluation of \mathcal{H}_T .

Method	\mathcal{H}_R	\mathcal{H}_T	Number of selected attributes	Run time (minutes)
<i>O1O2O3O4-NSGA-II</i>	0.0061	0.0115	10	38.03
<i>M6</i>	0.0079	0.0130	6	2.87
<i>M1</i>	0.0075	0.0141	28	5.08
<i>All attributes</i>	0.0112	0.0168	57	0.10
<i>M5</i>	0.0112	0.0168	57	0.09
<i>M2</i>	0.0089	0.0219	9	5.13
<i>M4</i>	0.0321	0.0394	1	29.44
<i>M3</i>	0.0321	0.0394	1	9.12

V. ANALYSIS OF RESULTS AND DISCUSSION

This section provides an analysis of the results presented above attending to the comparison of the multi-objective optimization models and the MOEAs that solve them, as well as the identification and adjustment of the final prediction

TABLE XI: Results of the best prediction model (obtained with *O1O2O3O4-NSGA-II* method) for the indoor temperature problem evaluated on *R* and *T* datasets.

Set	Metric	1-step ahead	2-steps ahead	3-steps ahead	4-steps ahead
<i>R</i>	RMSE	0.0079	0.0094	0.0108	0.0126
	MAE	0.0052	0.0064	0.0075	0.0089
	CC	0.9993	0.999	0.9987	0.9983
<i>T</i>	RMSE	0.0129	0.016	0.0194	0.0235
	MAE	0.0112	0.0139	0.0169	0.0204
	CC	0.9994	0.9991	0.9988	0.9983

model and its comparison with prediction models obtained with other FS techniques.

A. Analysis of the multi-objective optimization models

The experiments aimed at identifying the best multi-objective optimization model have been based on the evaluation of the prediction models obtained and on statistical tests to discard those models that have presented statistically significant differences. We have used the performance measure \mathcal{H} , proposed in this paper, to evaluate the prediction models on a dataset *E* (Algorithm 2, where *E* has been the dataset *R*, *V* or *T*) taking into account the RMSE, MAE and CC indicators in multi-step ahead predictions. For the statistical tests, the prediction models have been evaluated by means of 3 repetitions of 10-fold cross-validation, therefore using a total of 30 prediction models for each multi-objective optimization model. The following analysis can be derived:

- The combination of the evaluation of the surrogate model together with the correlation filter (multi-objective optimization model *O1O2O3*) produces the best prediction model (evaluated on the *T* dataset) in the air quality problem, while in the indoor temperature problem, the best resulting optimization model combines the evaluation of the surrogate model with the correlation and reliefF filters (multi-objective optimization model *O1O2O3O4*). This indicates that all the objective functions defined in this paper can be important in the multi-objective search for attribute subsets in the problems under study. Statistical tests have confirmed the superiority of these optimization models in their respective application problems.
- The *O1O2O3* optimization model for air quality forecasting has selected 16 attributes, and thus has removed 68 attributes from the original dataset (out of 84 attributes), as well as improving the performance of the LSTM neural network. The *O1O2O3O4* optimization model for indoor temperature forecasting has selected 10 attributes, removing 47 attributes from the original dataset (out of 57 attributes) and also improved the performance of the LSTM neural network.
- With the *O1O2* optimization model, which does not contain the correlation filter or the reliefF filter, a lower number of attributes is selected than in the rest of the optimization models, in both air quality and indoor temperature problems. However, the performance of the

prediction model obtained deteriorates with respect to other models that do contain some of these filters.

- The multi-objective optimization models *O1O2* and *O2O3O4* have the shortest run times, while the multi-objective optimization model *O1O2O3O4* containing all objectives has, as expected, the longest run time.

B. MOEAs comparison

In order to identify which MOEA has better performance, several outstanding MOEAs from the literature have been compared using optimality and diversity metrics used in multi-objective optimization. The multi-objective optimization models *O1O2O3*, for the air quality problem, and *O1O2O3O4*, for the indoor temperature problem, have been used in the comparison of MOEAs. Again, statistical tests have been performed to select, in this case, those MOEAs that have never lost in any of the statistical tests. For both multi-objective optimization models *O1O2O3* and *O1O2O3O4* in their respective application problems, the NSGA-II algorithm turned out to be the best. Figures 3 and 4 show the Pareto fronts found with *O1O2O3-NSGA-II* for air quality dataset and *O1O2O3O4-NSGA-II* for indoor temperature dataset, marking in red the solution selected by the decision-making process.

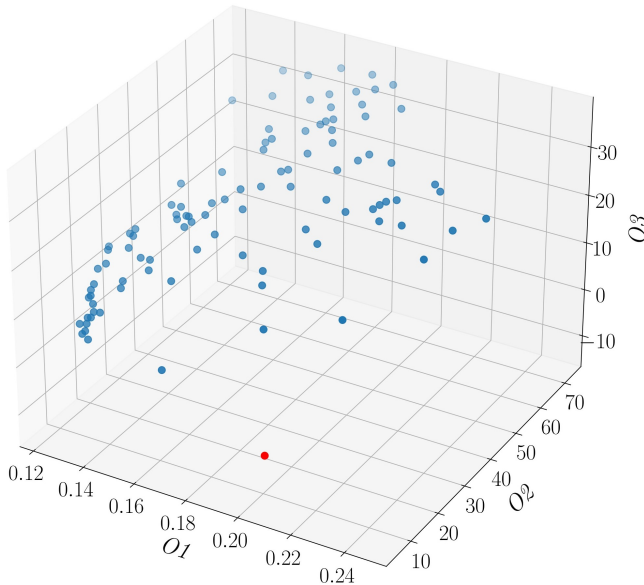


Fig. 3: Pareto front found with *O1O2O3-NSGA-II* (air quality dataset), best run, 100,000 evaluations. Red point represents the model with best average RMSE of 7-steps ahead predictions.

C. Analysis of the final prediction model

The performance measure \mathcal{H} proposed in this paper has been used finally to evaluate the prediction models on datasets R and T . We can highlight the following observations about the best prediction models found in this phase of experimentation:

- We can estimate the overfitting ratio of a model with the formula $\mathcal{H}_R/\mathcal{H}_T$. A rate greater than 1 indicates that the model is overfitting the training data, and a rate less than

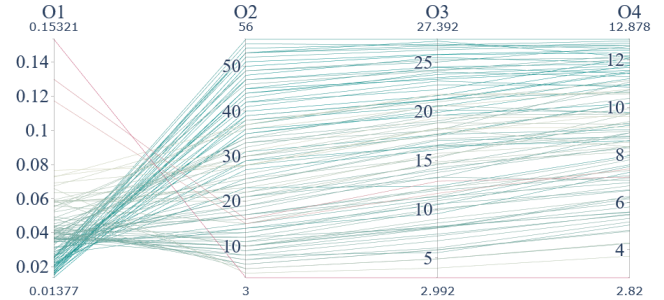


Fig. 4: Parallel coordinates plot of the Pareto front found with *O1O2O3O4-NSGA-II* (indoor temperature dataset), best run, 100,000 evaluations.

1 indicates that the model is under-fitting. The closer the rate is to 1, the better, since it implies that the results of training and testing are similar. The overfitting ratio of the prediction model found with *O1O2O3* and NSGA-II for the air quality problem is $\mathcal{H}_R/\mathcal{H}_T = 0.6217$, while the overfitting ratio of the prediction model found with *O1O2O3O4* and NSGA-II for the indoor temperature problem is $\mathcal{H}_R/\mathcal{H}_T = 0.5304$. This indicates that the proposed surrogate-assisted approach is not prone to overfitting, since there is no model training during the optimization phase, but rather an evaluation of the surrogate model on the validation dataset V .

- The prediction model found in this research with *O1O2O3* and NSGA-II for the air quality problem is, not only accurate but also stable, in the sense that the h -steps ahead predictions suffer only a small increase in the prediction error between one step and the next. The prediction model found with *O1O2O3O4* and NSGA-II for the indoor temperature problem is also accurate but suffers from less stability, mainly because the window size has been smaller in this case.

D. Comparison with other FS approaches

The following analysis can be derived from the comparison with the FS methods described in Section IV-F:

- None of the compared FS methods beats the *O1O2O3-NSGA-II* method nor the *O1O2O3O4-NSGA-II* method in their respective application problems.
- The *M1* and *M2* methods based on correlation and reliefF have performed better than the wrapper methods *M3* and *M4* based on linear regression and random forest. The *M1* and *M2* methods, in addition to using the correlation and reliefF filters respectively, build LSTM neural networks with the subsets of attributes, which can make a difference with respect to the *M3* and *M4* methods.
- The embedded method *M5* (based on the cancelOut layer) has not shown good performance, even after adjusting its hyperparameters. The *M5* method is the only FS method that has not selected any attributes.

- The embedded method *M6* (based on random forest) selects few attributes but performs worse than the best FS methods found in this paper.
- With regard to run times, the FS method *O1O2O3-NSGA-II* found in this work, which has evaluated 100,000 subsets of attributes (using the surrogate model) and has built a maximum of 100 LSTM neural networks (maximum number of non-dominated solutions in the last population), consumes less computational time than the *M4* wrapper method based on random forest, which builds and evaluates 100,000 prediction models. However, the *O1O2O3O4-NSGA-II* method has taken longer than the *M4* method due to the presence of 4 objective functions. The *M3* wrapper method based on linear regression, although it also builds 100,000 prediction models, consumes a similar run time to the *M1* and *M2* methods, since the linear regression learning algorithm is a really fast method. The embedded methods *M5* and *M6* based on cancelOut and random forest respectively are the least time consuming methods. The *M5* method consumes similar to the time required to train an LSTM neural network with all 84 attributes.

E. Interpretation of the selected attribute subsets

In this section we show the features that each model has selected for forecasting and discuss whether they really are important features from a domain perspective. This analysis strengthens the interpretability aspects of the models. Table XII shows the attributes selected by the *O1O2O3-NSGA-II* and *O1O2O3O4-NSGA-II* methods in the air quality and indoor temperature problems, respectively.

TABLE XII: Selected attributes in air quality an indoor temperature problems

Problem	Method	Selected Attributes
Air quality	<i>O1O2O3-NSGA-II</i>	<i>Lag_NO_1, Lag_NO_4, Lag_HR_7, Lag_NOx_1, Lag_NOx_2, Lag_NOx_3, Lag_NOx_4, Lag_NOx_6, Lag_NOx_7, Lag_DD_2, Lag_NO2_2, Lag_NO2_3, Lag_NO2_4, Lag_NO2_5, Lag_NO2_6, Lag_NO2_7</i>
Indoor temperature	<i>O1O2O3O4-NSGA-II</i>	<i>Lag_CO2_2, Lag_W_4, Lag_SLS_4, Lag_OT_1, Lag_OT_2, Lag_OT_4, Lag_IT_1, Lag_IT_2, Lag_IT_3, Lag_IT_4</i>

In the air quality problem, lagged variables of the attributes *NO*, *NO_x*, *DD* and *HR* were selected. The attributes *NO* and *NO_x* are closely related to the attribute *NO₂* in their chemical composition, and therefore there is a high (Pearson's) correlation between these attributes and the target attribute *NO₂*. Climate also influences *NO₂* concentrations, according to a study by the Leibniz Institute for Tropospheric Research (Germany) [66] commissioned by the State Office for Environment, Agriculture and Geology (LfULG). In this study was shown that wind speed and the height of the lowest air layer are the most important factors that determine how

much pollutants can accumulate locally. Moreover, it has long been known that weak winds can cause high concentrations of pollutants. The study also showed that high humidity can also reduce the concentration of *NO₂*, which could be due to the fact that the pollutants deposit more strongly on moist surfaces. The FS method has selected the first 6 lagged variables (out of 7), which may imply a low incidence in the prediction of the last day of the established time window.

In the indoor temperature problem, lagged variables of the attributes *CO₂*, *W*, *SLS* and *OT* have been selected. The presence of a high concentration of *CO₂* generated by the occupants of the room through breathing is a consequence of a poor ventilation system that leads to an increase in indoor temperature. On the other hand, it is evident that the external wind, the sunlight on the south facade (in areas of the northern hemisphere, as is the case) and the external temperature are factors that affect the indoor temperature of a room. With respect to the lagged variables of the target, all of them have been selected, which indicates a dependency of the entire time window on the model forecast.

F. Summary of analysis results

In summary, the results of the analysis are as follows:

- All the objectives *O1*, *O2*, *O3* and *O4* proposed in this research for multi-objective evolutionary feature selection have proven to be important in the identification of good forecasting models in the cases under study. Particularly, the combinations *O1O2O3* and *O1O2O3O4* have turned out to be the best in the problem of air quality and indoor temperature, respectively.
- The NSGA-II algorithm has shown the best performance, in terms of hypervolume, than the rest of the MOEAs in solving the multi-objective optimization problems proposed for FS in this paper.
- The prediction models found in this study with the proposed techniques present low overfitting rates and are stable with respect to multi-step ahead predictions.
- The proposed FS methods outperform other powerful filter-wrapper FS methods based on correlation, reliefF and LSTM, wrapper FS methods based on linear regression and random forest, and embedded FS methods based on cancelOut and random forest.
- The models found in this study for air quality prediction and indoor temperature prediction contain a relatively low number of attributes that are easy to interpret in the problem domain.

VI. CONCLUSIONS AND FUTURE WORK

In this paper, we have proposed new feature selection methods that are particularly appropriate for building prediction models that require a high computational training cost, such as models based on deep learning. The proposed methods search for subsets of attributes by means of a multi-objective evolutionary strategy and using an LSTM neural network which acts as a surrogate model to evaluate candidate subsets of attributes, together with correlation and reliefF filters and

cardinality minimization of attribute subsets. In the experimentation process, different optimization models of 2, 3 and 4 objectives have been tested, which combine the evaluation by means of the surrogate model with the correlation and reliefF filters. Different state-of-the-art multi-objective and many-objective evolutionary algorithms have also been tested, and the algorithm NSGA-II has shown the best performance. Two datasets has been used in the experiments, the first with air quality time series data in south-eastern Spain with data measured daily for 4 years, and the second with time series data related to the indoor temperature of a home domotic house. The data was divided into train, validation and test sets. To measure the performance of the prediction models, a performance measure multi-criteria metric has been proposed, which takes into account the average RMSE, MAE and CC in a multi-step ahead prediction horizon. In the air quality forecasting problem, the best FS method found was *OIO2O3-NSGA-II*, selecting 16 attributes out of a total of 84 and considerably improving the predictive capacity of the LSTM neural network with the complete set of attributes. In the indoor temperature forecasting problem, the best FS method found was *OIO2O3O4-NSGA-II*, which selected 10 attributes out of a total of 57, also improving the performance of the LSTM with all the attributes. They also show stability in the predictions without the presence of overfitting. The investigated FS methods *OIO2O3-NSGA-II* and *OIO2O3O4-NSGA-II* outperform, in their respective application scenario, FS methods of the wrapper, hybrid filter-wrapper, and embedded types. Finally, we have verified the importance of the selected attributes with the *OIO2O3-NSGA-II* and *OIO2O3O4-NSGA-II* methods in the expert context.

Future works will include the use of this novel approach with other deep learning algorithms such as GRU and CNN as well as other multi-step ahead forecasting strategies. In addition, predictions can be made for other harmful chemical compounds such as NO_X , $PM_{2.5}$, PM_{10} , etc. Besides, this approach can be applied to classification problems and imaging processing [67]. We are currently working on a new many-objective evolutionary algorithm based on decomposition, performance indicators and reference points sampled on the Pareto front, and on the development of a novel multi-surrogate assisted multi-objective evolutionary algorithm for feature selection method in order to improve the generalization error applied to time series forecasting problems. Finally, we are considering the application of the proposed technique in a spatio-temporal forecast scenario.

ACKNOWLEDGMENT

This work was partially funded by the CONFAINCE project (Ref: PID2021-122194OB-I00), supported by the Spanish Ministry of Science and Innovation and the Spanish Agency for Research, and the IMPACT-T2D project (PMP21/00092) supported by the Spanish Health Institute Carlos III (ISCIII).

REFERENCES

- [1] M. Kuhn and K. Johnson, *Feature Engineering and Selection: A Practical Approach for Predictive Models*, ser. Chapman & Hall/CRC Data Science Series. CRC Press, 2019.
- [2] K. Deb, *Multi-Objective Optimization Using Evolutionary Algorithms*. John Wiley & Sons, Inc., 2001.
- [3] S. Hochreiter and J. Schmidhuber, "Long Short-term Memory," *Neural computation*, vol. 9, pp. 1735–80, 12 1997.
- [4] M. Kowalska, M. Skrzypek, M. Kowalski, and J. Cyrus, "Effect of nox and no2 concentration increase in ambient air to daily bronchitis and asthma exacerbation, silesian voivodeship in poland," *International Journal of Environmental Research and Public Health*, vol. 17, no. 3, 2020.
- [5] R. W. Atkinson, B. K. Butland, H. R. Anderson, and R. L. Maynard, "Long-term concentrations of nitrogen dioxide and mortality: A meta-analysis of cohort studies," *Epidemiology (Cambridge, Mass.)*, vol. 29, no. 4, pp. 460–472, Jul 2018.
- [6] Kamath, U. and Liu, J., *Explainable Artificial Intelligence: An Introduction to Interpretable Machine Learning*. Springer International Publishing, 2021.
- [7] H. Liu and H. Motoda, *Feature Selection for Knowledge Discovery and Data Mining*. USA: Kluwer Academic Publishers, 1998.
- [8] Roberto Ruiz and Jesús S. Aguilar-Ruiz and José C. Riquelme and Norberto Díaz-Díaz, "Analysis of Feature Rankings for Classification," in *Lecture Notes in Computer Science*. Springer Berlin Heidelberg, 2005, pp. 362–372.
- [9] I. Kononenko, E. Šimec, and M. Robnik-Sikonja, "Overcoming the Myopia of Inductive Learning Algorithms with RELIEFF," *Applied Intelligence*, vol. 7, no. 1, pp. 39–55, Jan. 1997.
- [10] X. Zhang, M. Fan, D. Wang, P. Zhou, and D. Tao, "Top-k feature selection framework using robust 0–1 integer programming," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 7, pp. 3005–3019, 2021.
- [11] E. Schaffernicht, C. Möller, K. Debes, and H.-M. Gross, "Forward Feature Selection Using Residual Mutual Information," in *ESANN*, 01 2009, pp. 583–588.
- [12] I. Guyon, J. Weston, S. Barnhill, and V. Vapnik, "Gene selection for cancer classification using support vector machines," *Machine Learning*, vol. 46, pp. 389–422, 2004.
- [13] F. Jiménez, G. Sánchez, J. García, G. Sciacivco, and L. Miralles, "Multi-objective evolutionary feature selection for online sales forecasting," *Neurocomputing*, vol. 234, pp. 75–92, 2017.
- [14] M. Hall, "Correlation-Based Feature Selection for Machine Learning," *Department of Computer Science*, vol. 19, 06 2000.
- [15] M. Dash and H. Liu, "Consistency-Based Search in Feature Selection," *Artif. Intell.*, vol. 151, no. 1-2, pp. 155–176, Dec. 2003.
- [16] L. Yu and H. Liu, "Redundancy Based Feature Selection for Microarray Data," in *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ser. KDD '04. New York, NY, USA: Association for Computing Machinery, 2004, pp. 737–742.
- [17] S. Lei, "A Feature Selection Method Based on Information Gain and Genetic Algorithm," in *2012 International Conference on Computer Science and Electronics Engineering*, vol. 2, 2012, pp. 355–358.
- [18] K. Pearson, "X. On the criterion that a given system of deviations from the probable in the case of a correlated system of variables is such that it can be reasonably supposed to have arisen from random sampling," *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science*, vol. 50, no. 302, pp. 157–175, 1900.
- [19] R. Tibshirani, "Regression Shrinkage and Selection Via the Lasso," *Journal of the Royal Statistical Society, Series B*, vol. 58, pp. 267–288, 1994.
- [20] J. R. Quinlan, *C4.5: Programs for Machine Learning*. San Francisco, CA, USA: Morgan Kaufmann Publishers Inc., 1993.
- [21] L. Breiman, J. H. Friedman, R. A. Olshen, and C. J. Stone, *Classification and Regression Trees*. Chapman and Hall/CRC, 1984.
- [22] V. Bolón-Canedo and A. Alonso-Betanzos, "Recent Advances in Ensembles for Feature Selection," in *Intelligent Systems Reference Library*, 2018, pp. 13–37.
- [23] V. Borisov, J. Haug, and G. Kasneci, "Cancelout: A layer for feature selection in deep neural networks," in *Artificial Neural Networks and Machine Learning – ICANN 2019: Deep Learning*, I. V. Tetko, V. Kůrková, P. Karpov, and F. Theis, Eds. Cham: Springer International Publishing, 2019, pp. 72–83.
- [24] Y. Huang, W. Jin, Z. Yu, and B. Li, "Supervised feature selection through deep neural networks with pairwise connected structure," *Knowledge-Based Systems*, vol. 204, p. 106202, 2020.
- [25] L. Zhou, C. Zhang, M. F. Taha, X. Wei, Y. He, Z. Qiu, and Y. Liu, "Wheat Kernel Variety Identification Based on a Large Near-Infrared Spectral Dataset and a Novel Deep Learning-Based Feature Selection Method," *Frontiers in Plant Science*, vol. 11, p. 1682, 2020.

- [26] L. Cun, B. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. Hubbard, and L. D. Jackel, "Handwritten Digit Recognition with a Back-Propagation Network," in *Advances in Neural Information Processing Systems*. Morgan Kaufmann, 1990, pp. 396–404.
- [27] L. Wang, X. Yan, M.-L. Liu, K.-J. Song, X.-F. Sun, and W.-W. Pan, "Prediction of RNA-protein interactions by combining deep convolutional neural network with feature selection ensemble method," *Journal of Theoretical Biology*, vol. 461, pp. 230–238, 2019.
- [28] X. Yuan, L. Li, Y. A. Shardt, Y. Wang, and C. Yang, "Deep learning with spatiotemporal attention-based LSTM for industrial soft sensor model development," *IEEE Transactions on Industrial Electronics*, vol. 68, no. 5, pp. 4404–4414, 2020.
- [29] J. Wang, H. Zhang, J. Wang, Y. Pu, and N. R. Pal, "Feature Selection Using a Neural Network With Group Lasso Regularization and Controlled Redundancy," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 3, pp. 1110–1123, 2021.
- [30] W. Zheng, S. Chen, Z. Fu, F. Zhu, H. Yan, and J. Yang, "Feature Selection Boosted by Unselected Features," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 33, no. 9, pp. 4562–4574, 2022.
- [31] X. Wu, X. Xu, J. Liu, H. Wang, B. Hu, and F. Nie, "Supervised Feature Selection With Orthogonal Regression and Feature Weighting," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 5, pp. 1831–1838, 2021.
- [32] P. Zhou, P. Li, S. Zhao, and X. Wu, "Feature Interaction for Streaming Feature Selection," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 32, no. 10, pp. 4691–4702, 2021.
- [33] H. Ishibuchi, Q. Zhang, R. Cheng, K. Li, H. Li, H. Wang, and A. Zhou, Eds., *Evolutionary Multi-Criterion Optimization - 11th International Conference, EMO 2021, Shenzhen, China, March 28-31, 2021, Proceedings*, ser. Lecture Notes in Computer Science, vol. 12654. Springer, 2021.
- [34] Q. Al-Tashi, S. J. Abdulkadir, H. M. Rais, S. Mirjalili, H. Alhussian, M. G. Ragab, and A. Alqushaibi, "Binary Multi-Objective Grey Wolf Optimizer for Feature Selection in Classification," *IEEE Access*, vol. 8, pp. 106 247–106 263, 2020.
- [35] T. Niu, J. Wang, H. Lu, W. Yang, and P. Du, "Developing a deep learning framework with two-stage feature selection for multivariate financial time series forecasting," *Expert Systems with Applications*, vol. 148, p. 113237, 2020.
- [36] L. Shu, F. He, X. Hu, and H. Li, "A Novel Feature Selection with Many-Objective Optimization and Learning Mechanism," in *2021 IEEE 24th International Conference on Computer Supported Cooperative Work in Design (CSCWD)*, 2021, pp. 684–689.
- [37] K. Deb and H. Jain, "An evolutionary many-objective optimization algorithm using reference-point-based nondominated sorting approach, part i: Solving problems with box constraints," *IEEE Transactions on Evolutionary Computation*, vol. 18, no. 4, pp. 577–601, Aug 2014.
- [38] Y. Jin, "Surrogate-assisted evolutionary computation: Recent advances and future challenges," *Swarm and Evolutionary Computation*, vol. 1, no. 2, pp. 61–70, 2011.
- [39] Z. Lv, L. Wang, Z. Han, J. Zhao, and W. Wang, "Surrogate-assisted particle swarm optimization algorithm with pareto active learning for expensive multi-objective optimization," *IEEE/CAA Journal of Automatica Sinica*, vol. 6, no. 3, pp. 838–849, 2019.
- [40] T. H. Nguyen, D. Nong, and K. Paustian, "Surrogate-based multi-objective optimization of management options for agricultural landscapes using artificial neural networks," *Ecological Modelling*, vol. 400, pp. 1–13, 2019.
- [41] Z. Jiang, Y. Zhang, and J. Wang, "A multi-surrogate-assisted dual-layer ensemble feature selection algorithm," *Applied Soft Computing*, vol. 110, p. 107625, 2021.
- [42] H. Liu, Z. Duan, and C. Chen, "A hybrid multi-resolution multi-objective ensemble model and its application for forecasting of daily $PM_{2.5}$ concentrations," *Information Sciences*, vol. 516, pp. 266–292, 2020.
- [43] K. Deb, A. Pratab, S. Agarwal, and T. Meyarivan, "A fast and elitist multiobjective genetic algorithm: Nsga-ii," *IEEE Transactions on Evolutionary Computation*, vol. 6, no. 2, pp. 182–197, 2002.
- [44] S. Zhang, Y. Chen, W. Zhang, and R. Feng, "A novel ensemble deep learning model with dynamic error correction and multi-objective ensemble pruning for time series forecasting," *Information Sciences*, vol. 544, pp. 427–445, 2021.
- [45] J. Wang, H. Li, H. Yang, and Y. Wang, "Intelligent multivariable air-quality forecasting system based on feature selection and modified evolving interval type-2 quantum fuzzy neural network," *Environmental Pollution*, vol. 274, p. 116429, 2021.
- [46] P. Du, J. Wang, Y. Hao, T. Niu, and W. Yang, "A novel hybrid model based on multi-objective Harris hawks optimization algorithm for daily $PM_{2.5}$ and PM_{10} forecasting," *Applied Soft Computing*, vol. 96, p. 106620, 2020.
- [47] Fernando Jiménez and José Palma and Gracia Sánchez and David Marín and M.D. Francisco Palacios and M.D. Lucía López, "Feature selection based multivariate time series forecasting: An application to antibiotic resistance outbreaks prediction," *Artificial Intelligence in Medicine*, vol. 104, p. 101818, 2020.
- [48] Aurora González-Vidal and Fernando Jiménez and Antonio F. Gómez-Skarmeta, "A methodology for energy multivariate time series forecasting in smart buildings based on feature selection," *Energy and Buildings*, vol. 196, pp. 71–82, 2019.
- [49] R. Espinosa, J. Palma, F. Jiménez, J. Kaminska, G. Sciavicco, and E. Lucena-Sánchez, "A time series forecasting based multi-criteria methodology for air quality prediction," *Applied Soft Computing*, vol. 113, p. 107850, 2021.
- [50] Raquel Espinosa and Fernando Jiménez and José Palma, "Multi-objective evolutionary spatio-temporal forecasting of air pollution," *Future Generation Computer Systems*, vol. 136, pp. 15–33, 2022.
- [51] Xue, Bing and Zhang, Mengjie and Browne, Will N. and Yao, Xin, "A survey on evolutionary computation approaches to feature selection," *IEEE Transactions on Evolutionary Computation*, vol. 20, no. 4, pp. 606–626, 2016.
- [52] L. J. Eshelman, "The CHC Adaptive Search Algorithm : How to Have Safe Search When Engaging in Nontraditional Genetic Recombination," *Foundations of Genetic Algorithms*, vol. 1, pp. 265–283, 1991.
- [53] L. Davis, Ed., *Handbook of Genetic Algorithms*. Van Nostrand Reinhold, 1991.
- [54] G. Bontempi, S. B. Taieb, and Y. Borgne, "Machine learning strategies for time series forecasting," in *eBISS*, 2013, pp. 62–77.
- [55] Taieb, Souhaib Ben and Atiya, Amir F., "A Bias and Variance Analysis for Multistep-Ahead Time Series Forecasting," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 27, no. 1, pp. 62–76, 2016.
- [56] Ruofeng Wen and Kari Torkkola and Balakrishnan (Murali) Narayanaswamy and Dhruv Madeka, "A multi-horizon quantile recurrent forecaster," in *NeurIPS 2017*, 2017.
- [57] J. Brownlee, "Time Series Forecasting as Supervised Learning," Dec. 2020. [Online]. Available: <https://machinelearningmastery.com/time-series-forecasting-supervised-learning/>
- [58] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *Icml*, 2010.
- [59] Q. Zhang and H. Li, "MOEA/D: A Multiobjective Evolutionary Algorithm Based on Decomposition," *IEEE Transactions on Evolutionary Computation*, vol. 11, no. 6, pp. 712–731, 2007.
- [60] E. Zitzler and S. Künzli, "Indicator-based selection in multiobjective search," in *Parallel Problem Solving from Nature - PPSN VIII*, X. Yao, E. K. Burke, J. A. Lozano, J. Smith, J. J. Merelo-Guervós, J. A. Bullinaria, J. E. Rowe, P. Tiño, A. Kabán, and H.-P. Schwefel, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2004, pp. 832–842.
- [61] Z. Fan, W. Li, X. Cai, H. Huang, Y. Fang, Y. You, J. Mo, C. Wei, and E. Goodman, "An Improved Epsilon Constraint-handling Method in MOEA/D for CMOPs with Large Infeasible Regions," 2017.
- [62] E. Zitzler, M. Laumanns, and L. Thiele, "Spea2: Improving the strength pareto evolutionary algorithm," *TIK-report*, vol. 103, 2001.
- [63] J. B. Kollat and P. M. Reed, "The value of online adaptive search: A performance comparison of nsgaii, ϵ -nsgaii and emoea," in *Evolutionary Multi-Criterion Optimization*, C. A. Coello Coello, A. Hernández Aguirre, and E. Zitzler, Eds. Berlin, Heidelberg: Springer Berlin Heidelberg, 2005, pp. 386–398.
- [64] Dua, Dheeru and Graff, Casey, "UCI Machine Learning Repository," 2017. [Online]. Available: <http://archive.ics.uci.edu/ml>
- [65] F. Zamora-Martínez and P. Romeu and P. Botella-Rocamora and J. Pardo, "On-line learning of indoor temperature forecasting models towards energy efficiency," *Energy and Buildings*, vol. 83, pp. 162–172, 2014.
- [66] Dominik van Pinxteren and Sebastian Düsing and Alfred Wiedensohler and Hartmut Herrmann, "Meteorological influences on nitrogen dioxide: Influence of weather conditions and weathering on nitrogen dioxide concentrations in outdoor air 2015 to 2018," *Series of publications of the LfJULG*, vol. 2/2020, 2020.
- [67] A. Glowacz, "Fault diagnosis of electric impact drills using thermal imaging," *Measurement*, vol. 171, p. 108815, 2021.